エントロピ正則された強化学習を用いた模倣学習

Imitation learning based on entropy-regularized reinforcement learning

内部 英治 *1 Eiji Uchibe

*1国際電気通信基礎技術研究所

ATR Computational Neuroscience Labs.

This paper proposes Entropy-Regularized Imitation Learning (ERIL) that is given by a combination of forward and inverse reinforcement learning. ERIL utilizes the soft Bellman optimality equation in which the reward function is augmented by the entropy of the learning policy and the Kullback-Leibler (KL) divergence between the learning and the baseline policies. We show that inverse RL is interpreted as estimating the log-ratio between two policies and the log-ratio is efficiently solved by binary logistic regression. Forward RL is given by a variant of Dynamic Policy Programming and our algorithm is interpreted as minimization of the KL divergence between the learning policy and the estimated expert policy. Experimental results on the MuJoCo-simulated environments show that ERIL is more sample efficient than the previous methods such as GAIL and AIRL because the forward RL step of ERIL is off-policy.

1. はじめに

強化学習は最適方策(制御則)を試行錯誤によって獲得する ための手法で,知能ロボットに自律的な学習能力を与えるため に必要である.またヒトや動物の意思決定過程の計算モデルと して計算論的神経科学の分野で注目されている.一般に強化学 習では各状態で実行した行動に対する即時評価である報酬を必 要とする.しかし所望の最適方策を達成する報酬を設計するの は簡単ではない.たとえばタスクを達成した場合にのみ正の報 酬値を与え,それ以外は0であるような非常にスパースな報酬 関数を用いた場合,原理的には最適方策を学習できるが必要と するデータ数や学習時間は非常に膨大なものとなる.学習効率 を改善するために複雑な報酬関数を設計することも考えられる が,多くの場合は意図しない最適方策を学習することになる.

```
この問題に対処するための方法として逆強化学習がある. 逆
強化学習ではタスクを達成しているエキスパートの存在を仮
定し、エキスパートが用いている報酬関数をエキスパートか
らのデータから推定する.近年,強化学習と逆強化学習の組
み合わせが敵対的生成ネットワーク (Generative Adversarial
Networks; GAN) として定式化できることが示された [3, 4, 6].
逆強化学習は GAN における識別器に対応し、エキスパートか
らのデータとロボット自身が生成したデータを区別する.強化
学習は GAN における生成器に対応し、逆強化学習によって推
定された報酬の期待積算を最大にするように新しい方策を学
習する. Generative Adversarial Imitation Learning (GAIL)
[6] Adversarial Inverse Reinforcement Learning (AIRL)
[4] は逆強化学習と強化学習を繰り返すことで単純な模倣学習で
ある Behavior Cloning (BC) よりもエキスパート方策を精度
良く復元できることを示した. GAIL や AIRL はエキスパート
からのサンプル数が少ない場合でも有効であるが、その一方で
生成器の改善に必要な環境とのインタラクションに関するデー
タ効率は低いことが知られている.一つの理由として GAIL や
AIRL は方策オン型の強化学習法である Trust Region Policy
Optimization (TRPO) [9] を用いていることがあげられる.重
```

連絡先: 内部英治,国際電気通信基礎技術研究所 脳情報研究所 ブレインロボットインタフェース研究室,〒 619-0288 京 都府相楽郡精華町光台二丁目 2 番地 2, uchibe@atr.jp 点サンプリング等の技術を用いない限り過去に学習方策が生成 したデータを再利用できないため方策の改善には現在の学習方 策が多くのデータを生成する必要がある.

そこで強化学習時のデータ効率を改善するために、本研究で はエントロピ正則された模倣学習 (Entropy-Regularized Imitation Learning; ERIL) を提案する.報酬関数は学習方策の エントロピと学習方策とベースライン方策の Kullback-Leibler (KL) ダイバージェンスによって正則化される. この定式化の もとで得られるソフトベルマン最適方程式を用いて、逆強化 学習を密度比推定問題として定式化する.密度比推定はエキ スパート方策からのデータと学習者自身が生成したデータを 区別するロジスティック回帰によって効率的に解くことができ る. 識別器は報酬, 状態価値関数, 学習者の方策によって記述 され,従来法で用いられた識別器 [4, 12] を一般化したもので ある. 逆強化学習の結果エキスパート方策が推定されるため, KL ダイバージェンスによって新たな学習方策を更新する、こ れは方策オフ型の強化学習である Soft Actor-Critic [5] で用い られるものと同一で, 強化学習時におけるデータ効率を期待で きる.

OpenAI gym [2] で提供されているロボット制御課題を用い て ERIL と従来手法を比較する.シミュレーション結果より, エキスパートからのデータ数に関しては従来法と同程度の効率 であるが,環境とのインタラクションに関するデータ効率は従 来法よりも改善されたことを示す.

2. エントロピ正則化された強化学習

離散時間の無限期間マルコフ決定過程 (Markov Decision Process; MDP) を考える. MDP は $(\mathcal{X}, \mathcal{U}, p_T, r, \gamma, p_0)$ の組に よって定義される. ここで \mathcal{X}, \mathcal{U} はそれぞれ状態空間,行動空 間である. $p_T(\mathbf{x}' \mid \mathbf{x}, \mathbf{u})$ は状態 $\mathbf{x} \in \mathcal{X}$ で行動 $\mathbf{u} \in \mathcal{U}$ を実行し た時に状態 $\mathbf{x}' \in \mathcal{X}$ に遷移する確率で,モデルフリー強化学習 の枠組みでは未知である. $\hat{r}(\mathbf{x}, \mathbf{u})$ は状態 \mathbf{x} ,行動 \mathbf{u} に対して 与えられる即時報酬である. $\gamma \in (0, 1)$ は割引率, $p_0(\mathbf{x})$ は初期 状態分布である. 状態 \mathbf{x} で行動 \mathbf{u} を選択する確率を $\pi(\mathbf{u} \mid \mathbf{x})$ とする. 強化学習の目的は期待割引積算報酬を最大にする方策 を求めることである. 近年,報酬関数をエントロピによって正則化されたクラスの MDP が注目を集めている [1, 5, 7]. 具体的には報酬関数が

$$\tilde{r}(\boldsymbol{x}, \boldsymbol{u}) \triangleq r(\boldsymbol{x}, \boldsymbol{u}) + \kappa \mathcal{H}(\pi(\cdot \mid \boldsymbol{x})) - \eta \mathrm{KL}(\pi(\cdot \mid \boldsymbol{x}) \parallel b(\cdot \mid \boldsymbol{x}))$$
(1)

のように正則化される. ここで r(x, u) は通常の意味での即時 報酬関数, $\mathcal{H}(\pi(\cdot \mid x))$ は方策 π のエントロピ, $\mathrm{KL}(\pi(\cdot \mid x) \mid b(\cdot \mid x))$ は π とベースライン方策 $b(u \mid x)$ の間の Kullback-Leibler (KL) ダイバージェンス, κ, η は実験者の定めるメタ パラメータである. このとき以下のベルマン最適性方程式

$$V(\boldsymbol{x}) = \max_{\pi} \mathbb{E}_{\boldsymbol{u} \sim \pi(\cdot \mid \boldsymbol{x})} \left[r(\boldsymbol{x}, \boldsymbol{u}) - \kappa \ln \pi(\boldsymbol{u} \mid \boldsymbol{x}) - \eta \ln \frac{\pi(\boldsymbol{u} \mid \boldsymbol{x})}{b(\boldsymbol{u} \mid \boldsymbol{x})} + \gamma \mathbb{E}_{\boldsymbol{x}' \sim p_T(\cdot \mid \boldsymbol{x}, \boldsymbol{u})} \left[V(\boldsymbol{x}') \right] \right]$$
(2)

が成り立つ. ここで $V(\mathbf{x})$ は状態価値関数と呼ばれる. エン トロピの役割は最適方策が決定論的になることを防ぎ,探索 を促進する. KL ダイバージェンスの役割はベースライン方策 $b(\mathbf{u} \mid \mathbf{x})$ からあまり逸脱しないように方策改善ステップを保守 的にする. エントロピ正則化の利点は式 (2) の右辺の方策に関 する最大化が Lagrange の未定乗数法によって解析的に求めら れることである. 結果として状態価値関数 $V(\mathbf{x})$ と状態行動価 値関数 (\mathbf{x}, \mathbf{u}) に対して以下の関係式

$$V(\boldsymbol{x}) = \beta^{-1} \ln \int \exp\left(\beta Q(\boldsymbol{x}, \boldsymbol{u})\right) d\boldsymbol{u}, \qquad (3)$$

$$Q(\boldsymbol{x}, \boldsymbol{u}) = r(\boldsymbol{x}, \boldsymbol{u}) + \eta \ln b(\boldsymbol{u} \mid \boldsymbol{x}) + \gamma \mathbb{E}_{\boldsymbol{x}' \sim p_T(\cdot \mid \boldsymbol{x}, \boldsymbol{u})} \left[V(\boldsymbol{x}') \right].$$
(4)

が成り立つ. ここで

$$\beta \triangleq \frac{1}{\kappa + \eta}$$

と定義している. 行動が離散の場合,式 (3)の右辺は log-sumexp 関数として知られる max 演算子を滑らかにしたものであ る. 最適方策は

$$\pi(\boldsymbol{u} \mid \boldsymbol{x}) = \exp\left[\beta\left(Q(\boldsymbol{x}, \boldsymbol{u}) - V(\boldsymbol{x})\right)\right], \quad (5)$$

と与えられる. $\eta = 0$ のとき Soft Q-learning や Soft Actor-Critic [5] が, $\kappa = 0$ のとき Dynamic Policy Programming (DPP) [1] が導出される.

3. エントロピ正則化された模倣学習

エントロピ正則化された強化学習は報酬関数 r(x, u) が与え られれば最適方策を求めることができる.本節では報酬関数の かわりに最適方策から生成された状態行動データから最適方策 を求める模倣学習を提案する.

3.1 密度比推定による逆強化学習

式 (1) において最適方策を $\pi(\boldsymbol{u} \mid \boldsymbol{x}) \equiv \pi^{E}(\boldsymbol{u} \mid \boldsymbol{x}), \quad \forall - \boldsymbol{\lambda}$ ライン方策を $b(\boldsymbol{u} \mid \boldsymbol{x}) \equiv \pi^{G}(\boldsymbol{u} \mid \boldsymbol{x})$ とする. このとき式 (4), (5) より以下の関係式

$$\frac{1}{\beta} \ln \frac{\pi^{E}(\boldsymbol{u} \mid \boldsymbol{x})}{\pi^{G}(\boldsymbol{u} \mid \boldsymbol{x})} = r(\boldsymbol{x}, \boldsymbol{u}) - \kappa \ln \pi^{G}(\boldsymbol{u} \mid \boldsymbol{x}) + \gamma \mathbb{E}_{\boldsymbol{x}' \sim p_{T}}(\cdot \mid \boldsymbol{x}, \boldsymbol{u}) \left[V(\boldsymbol{x}') \right] - V(\boldsymbol{x}) \quad (6)$$

を得る. さらに $\pi^{E}(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}')$ を

$$\pi^{E}(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}') \triangleq p_{T}(\boldsymbol{x}' \mid \boldsymbol{x}, \boldsymbol{u}) \pi^{E}(\boldsymbol{u} \mid \boldsymbol{x}) \pi^{E}(\boldsymbol{x})$$

と提議し, $\pi^G(\mathbf{x}, \mathbf{u}, \mathbf{x}')$ も同様に定義する. このとき式 (6) の 左辺は

$$\ln \frac{\pi^{E}(\boldsymbol{u} \mid \boldsymbol{x})}{\pi_{k}^{G}(\boldsymbol{u} \mid \boldsymbol{x})} = \ln \frac{\pi^{E}(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}')}{\pi_{k}^{G}(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}')} - \ln \frac{\pi^{E}(\boldsymbol{x})}{\pi_{k}^{G}(\boldsymbol{x})}$$
(7)

と書き換えることができる.式 (7)の右辺は二つの密度比の差 であり、それぞれ π^E, π^G からのサンプルがあれば推定できる [10].たとえば $\pi^E(\mathbf{u} \mid \mathbf{x})$ から収集された状態行動遷移を

$$\mathcal{D}^{E} = \{ (\boldsymbol{x}_{i}, \boldsymbol{u}_{i}, \boldsymbol{x}'_{i}) \}_{i=1}^{N^{E}},$$
$$\boldsymbol{u}_{i} \sim \pi^{E}(\cdot \mid \boldsymbol{x}_{i}), \quad \boldsymbol{x}'_{i} \sim p_{T}(\cdot \mid \boldsymbol{x}_{i}, \boldsymbol{u}_{i})$$

とし、同様に π^{G} からのサンプル \mathcal{D}^{G} を収集する. 密度比推定 問題を通して報酬や状態価値関数を推定することが逆強化学習 に対応する.

密度比を推定する様々な手法が提案されている [10] が, ここで はロジスティック回帰を用いる.そのために二値変数 $L \in \{0,1\}$ を導入し, L = 1 のとき π^{E} からサンプル, L = 0 のとき π^{G} からサンプルする.すなわち

$$\pi^{G}(\boldsymbol{x}) = \Pr(\boldsymbol{x} \mid L=0), \quad \pi^{E}(\boldsymbol{x}) = \Pr(\boldsymbol{x} \mid L=1)$$

とする. このとき

$$\ln \frac{\pi^{E}(\boldsymbol{x})}{\pi^{G}(\boldsymbol{x})} = \ln \frac{D(\boldsymbol{x})}{1 - D(\boldsymbol{x})} - \ln \frac{\Pr(L=1)}{\Pr(L=0)}$$

と書き換えることができる. ここで $D(\mathbf{x})$ は状態 \mathbf{x} が π^E から のサンプルであるかどうかを識別する識別器 $D(\mathbf{x}) \triangleq \Pr(L = 1 | \mathbf{x})$ である. 同様に $\pi^E(\mathbf{x}, \mathbf{u}, \mathbf{x}')/\pi_k^G(\mathbf{x}, \mathbf{u}, \mathbf{x}')$ も識別器 $D(\mathbf{x}, \mathbf{u}, \mathbf{x}') \triangleq \Pr(L = 1 | \mathbf{x}, \mathbf{u}, \mathbf{x}')$ によって記述できる. 以上 より式 (6) を

$$\frac{1}{\beta} \ln \frac{D(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}')}{1 - D(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}')} = \frac{1}{\beta} \ln \frac{D(\boldsymbol{x})}{1 - D(\boldsymbol{x})} + r(\boldsymbol{x}, \boldsymbol{u}) - \kappa \ln \pi^{G}(\boldsymbol{u} \mid \boldsymbol{x}) + \gamma V(\boldsymbol{x}') - V(\boldsymbol{x}) \quad (8)$$

と書き換える.対数密度比 $\ln D(\boldsymbol{x})/(1 - D(\boldsymbol{x}))$ を $g(\boldsymbol{x})$ とおく. このとき識別器 $D(\boldsymbol{x})$ は

$$D(\boldsymbol{x}) = \frac{1}{1 + \exp(-g(\boldsymbol{x}))}$$

となり、通常の二値分類問題として g(x) を推定できる. 同様 にして D(x, u, x') を

$$D(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}') = \frac{\exp(\beta f(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}'))}{\exp(\beta f(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}')) + \exp(\beta \kappa \ln \pi^G(\boldsymbol{u} \mid \boldsymbol{x}))}$$
(9)

と構築する. ここで

$$f(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}') \triangleq r(\boldsymbol{x}, \boldsymbol{u}) + \frac{1}{\beta}g(\boldsymbol{x}) + \gamma V(\boldsymbol{x}') - V(\boldsymbol{x})$$

と定義している.式 (9)の識別器は従来の逆強化学習で用いられている識別器の拡張である.たとえば $g(\mathbf{x}) = 0$ かつ $\eta = 0$

とすれば AIRL [4] の識別器が, $\kappa = 0$ とすれば LogReg-IRL [12] の識別器と一致する.

以上をまとめると逆強化学習は次の二つのステップで構成 される.まず識別器 D(x) を学習することで g(x) を推定する. 次に識別器 D(x, u, x') を学習して r(x, u) と V(x) を推定する. この際 g(x) と $\pi^{G}(u \mid x)$ は固定する. D(x, u, x') は対 数尤度

$$J(D) = \mathbb{E}_{(\boldsymbol{x},\boldsymbol{u},\boldsymbol{x}')\sim\mathcal{D}^{E}} \left[\ln D(\boldsymbol{x},\boldsymbol{u},\boldsymbol{x}')\right] \\ + \mathbb{E}_{(\boldsymbol{x},\boldsymbol{u},\boldsymbol{x}')\sim\mathcal{D}^{G}} \left[\ln \left(1 - D(\boldsymbol{x},\boldsymbol{u},\boldsymbol{x}')\right)\right]$$

を最大にすることで得られる.ここで $(x, u, x') \sim D$ はDからのサンプルの期待値を意味する.

3.2 ベースライン方策の更新

密度比推定による逆強化学習の結果,エキスパート方策 π^E の推定値

$$\frac{1}{\beta} \ln \hat{\pi}^{E}(\boldsymbol{u} \mid \boldsymbol{x}) \approx r(\boldsymbol{x}, \boldsymbol{u}) + \eta \ln \pi^{G}(\boldsymbol{u} \mid \boldsymbol{x}) + \gamma V(\boldsymbol{x}') - V(\boldsymbol{x})$$

が得られる. これを用いて新しいベースライン方策
 $\pi^G_{\rm new}$ を KL ダイバージェンス

$$J(\pi) = \mathbb{E}\left[\mathrm{KL}(\pi(\cdot \mid \boldsymbol{x}) \parallel \hat{\pi}^{E}(\cdot \mid \boldsymbol{x}))\right]$$
(10)

を最小にするように求める.

ベースライン方策の更新は Soft Actor-Critic [5] で用いら れるものと同一である. Soft Actor-Critic では最適方策は以 下の KL ダイバージェンス

$$J(\pi) = \mathbb{E}\Big[\mathrm{KL}\Big(\pi(\cdot \mid \boldsymbol{x}) \parallel \exp\{\beta \left(Q(\boldsymbol{x}, \cdot) - V(\boldsymbol{x})\right)\}\Big)\Big],$$

を最小にすることで計算される. Soft Actor-Critic では報酬 r(x, u)が与えられたとき式 (3), (4)を用いてV(x)とQ(x, u)を計算したのちに最適方策を推定するのに対し,提案手法では エキスパートからのデータ \mathcal{D}^E が与えられたときに密度比推 定によって最適方策を推定することになる.よって式 (10)に よる方策改善が強化学習ステップに相当する.

ただし式 (10) による更新では \mathcal{D}^{E} と現在の学習方策 π^{G} に よって生成されたデータ \mathcal{D}^{G} だけから推定されたエキスパー ト方策を用いているため,過去の学習方策が用いたデータを 利用していない. そこで逆強化学習によって推定された報酬と 対数密度を固定し,状態価値関数と状態行動価値関数を Soft Actor-Critic にしたがって更新する.式 (5) より,状態価値関 数を更新するための目的関数は

$$J(V) = \mathbb{E}_{\boldsymbol{x} \sim \mathcal{D}} \left[\frac{1}{2} \left(V(\boldsymbol{x}) - \tilde{V}(\boldsymbol{x}) \right)^2 \right],$$

と与えられる. ここで

$$\tilde{V}(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{u} \sim \pi^G} \left[Q(\boldsymbol{x}, \cdot) - \beta^{-1} \ln \pi^G(\cdot \mid \boldsymbol{x}) \right]$$

である.状態行動価値関数を学習する目的関数は

$$J(Q) = \mathbb{E}_{(\boldsymbol{x},\boldsymbol{u},\boldsymbol{x}')\sim\mathcal{D}}\left[\frac{1}{2}\left\{Q(\boldsymbol{x},\boldsymbol{u}) - \tilde{Q}(\boldsymbol{x},\boldsymbol{u},\boldsymbol{x}')\right\}^2\right],$$

となる. ここで

$$\tilde{Q}(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{x}') = r(\boldsymbol{x}, \boldsymbol{u}) + \eta^{-1} \ln \pi^{G}(\boldsymbol{u} \mid \boldsymbol{x}) + \gamma V(\boldsymbol{x}')$$

とする. D はエキスパートデータも含めたすべてのデータを まとめたものである.



図 1: エキスパートからのデータ数についての性能比較. 横軸 はエキスパート方策 π^E から与えられた軌跡の個数, 縦軸は本 来の報酬のもとで評価した場合の性能を示す.

4. シミュレーション

提案手法の有効性を検証するために、OpenAI gym [2] で提 供されている Half-Cheetah, Hopper, Walker, Ant という 4 種類のロボット制御課題に適用する. これらは物理エンジン MuJoCo [11] を用いている. これらのタスクの目的はできる だけ速く移動することである.まず本来設定されている報酬関 数をもとに最適方策 π^E を TRPO によって学習し、そこから 得られるデータを D^E として用いる. 関数近似として用いる ニューラルネットワークの構造は従来研究 [4,6] を参考に構築 した. 対数密度 $g(\mathbf{x})$,報酬 $r(\mathbf{x})$,状態価値関数 $V(\mathbf{x})$,行動 価値関数 Q(x, u) は 2 層のニューラルネットワークを用い,活 性度関数は ReLU, ユニット数はそれぞれ 400, 300 とした *1. また方策 $\pi^{G}(\boldsymbol{u} \mid \boldsymbol{x})$ はガウス分布によって構成し、その平均 値を同じ構成のニューラルネットワークで表現した. 初期ベー スライン方策は BC によって初期化する. 1 エポックあたり 学習方策 π^G によって生成される軌跡は 100 とし, 各軌跡は 50 個の状態行動遷移対を含むとする. 正則化のパラメータは $\kappa = 1, \eta = 10$ とし、割引率は $\gamma = 0.99$ とした.

まず D^E からのサンプル数を変えることで逆強化学習のデー タ効率を評価する. Ho and Ermon [6] に従い, 一つの軌跡が 50 個の状態行動遷移対 (x, u, x') を含む, つまり1エピソード あたりのステップ数が 50 である. 逆強化学習は不良設定問題 であるため、各手法で推定された報酬は直接比較できない. そ こで最終的に獲得された学習方策をシミュレータで提供される 本来の報酬のもとで評価する.図1に ERIL と GAIL, AIRL. BC を比較した結果を示す. ERIL, GAIL, AIRL ともに BC よりも高い制御性能を示す方策を獲得している.一方で3種類 の模倣学習の間にはそれほど違いは見られなかった. Fu et al. [4] は GAIL と AIRL で実質的な差はなかったことを報告して おり、本報告の結果と一致する. また前述したように対数密度 $g(\mathbf{x})$ とメタパラメータ κ, η の値によって ERIL と AIRL の識 別器は一致する.実際にこの課題では学習方策をBCによって 初期化しているため、ほとんどの状態 x において $g(x) \approx 0$ と なっていた.

次にエキスパート方策 π^E からの軌跡の数は 25 と固定し,

^{*1} Fu et al. [4] が指摘したように逆強化学習では報酬を状態と行動 の関数として近似した場合は表現能力の高さから過剰適合しやすい. そのため報酬関数は状態の身に依存する関数として表現した.



図 2: 学習方策からのデータ数についての性能比較. 横軸は学 習方策 π^G から各エポックで与えられた軌跡の個数, 縦軸は本 来の報酬のもとで評価した場合の性能を示す.

学習方策が生成する軌跡の数を変更した場合の性能を比較した. 図 2 に結果を示す. 軌跡の数が少ない場合, ERIL は GAIL, AIRL, BC よりも高い性能を達成する学習方策を獲得した. こ の理由としては ERIL は方策オフ型の強化学習を用いている のに対し, GAIL と AIRL は方策オン型である TRPO を用い ているためであると考えられる.

5. おわりに

本稿ではエントロピ正則化強化学習に基づいた模倣学習 ERIL を提案した.GAN-GCL,GAIL,AIRL などの従来法と異な り,提案手法はソフトベルマン最適方程式に基づき導出されて いる.そのため強化学習に相当するベースラインの更新ステッ プは方策オフ型である.シミュレーション結果より ERIL は 従来法よりもサンプル効率が良いことが示された.

ERIL には注意深く調整すべきメタパラメータがあり,強化 学習の場合には式 (1) で用いている正則化のパラメータ κ, η は学習過程の安定性や関数近似誤差に影響を与えることが知ら れている [7]. 逆強化学習時においても識別器の性能に影響を 与えるが,これらのパラメータをどのように学習中に調整する かは今後の課題である.

強化学習ステップにおける ERIL と AIRL の性能差は主に Soft Actor-Critic と TRPO の違いによって説明できると考え られる. Neu et al. は DPP, TRPO など様々な強化学習ア ルゴリズムの性質を調査した [8] が, Soft Actor-Critic や逆強 化学習を組み合わせた場合の性質はいまだ明らかではない. こ の研究を拡張することで ERIL と AIRL の差を考察すること も今後の課題である.

謝辞 本研究の成果は,国立研究開発法人新エネルギー・産業 技術総合開発機構 (NEDO) の委託業務および JST,未来社会 創造事業,JPMJMI18B8 の支援を受けたものである.

参考文献

 M. G. Azar, V. Gómez, and H. J. Kappen. Dynamic policy programming. *Journal of Machine Learning Re*search, 13:3207–3245, 2012.

- [2] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. OpenAI Gym. arXiv preprint, 2016.
- [3] C. Finn, P. Christiano, P. Abbeel, and S. Levine. A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models. In *NIPS 2016 Workshop on Adversarial Training*, 2016.
- [4] J. Fu, K. Luo, and S. Levine. Learning robust rewards with adversarial inverse reinforcement learning. In Proc. of the 6th International Conference on Learning Representations, 2018.
- [5] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In Proc. of the 35th International Conference on Machine Learning, pages 1856–1865, 2018.
- [6] J. Ho and S. Ermon. Generative adversarial imitation learning. In Advances in Neural Information Processing Systems 29. 2016.
- [7] T. Kozuno, E. Uchibe, and K. Doya. Theoretical analysis of efficiency and robustness of softmax and gapincreasing operators in reinforcement learning. In Proc. of the 22nd International Conference on Artificial Intelligence and Statistics, 2019.
- [8] G. Neu, V. Gómez, and A. Jonsson. A unified view of entropy-regularized markov decision processes. arXiv preprint, 2017.
- [9] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In Proc. of the 32nd International Conference on Machine Learning, pages 1889–1897, 2015.
- [10] M. Sugiyama, T. Suzuki, and T. Kanamori. *Density ratio estimation in machine learning*. Cambridge University Press, 2012.
- [11] E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In Proc. of IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 5026–5033, 2012.
- [12] E. Uchibe. Model-free deep inverse reinforcement learning by logistic regression. *Neural Processing Letters*, 47(3):891–905, 2018.