

(19) 日本国特許庁(JP)

(12) 特許公報(B2)

(11) 特許番号

特許第4586129号
(P4586129)

(45) 発行日 平成22年11月24日(2010.11.24)

(24) 登録日 平成22年9月17日(2010.9.17)

(51) Int.Cl. F I
G05B 13/02 (2006.01) G05B 13/02 J
 G05B 13/02 L

請求項の数 9 (全 35 頁)

<p>(21) 出願番号 特願2008-77671 (P2008-77671) (22) 出願日 平成20年3月25日 (2008.3.25) (65) 公開番号 特開2009-230645 (P2009-230645A) (43) 公開日 平成21年10月8日 (2009.10.8) 審査請求日 平成21年2月12日 (2009.2.12)</p> <p>特許法第30条第1項適用 平成19年9月25日 http://isw3.naist.jp/IS/TechReport/を通じて発表、平成19年9月26日 http://library.naist.jp/library/tr/index.htmlを通じて発表、平成20年1月30日 平成19年度 国立大学法人奈良先端科学技術大学院大学 情報科学研究科 博士学位論文公聴会に発表</p>	<p>(73) 特許権者 506086797 独立行政法人沖縄科学技術研究基盤整備機構 沖縄県国頭郡恩納村字谷茶1919-1</p> <p>(74) 代理人 100064746 弁理士 深見 久郎</p> <p>(74) 代理人 100085132 弁理士 森田 俊雄</p> <p>(74) 代理人 100083703 弁理士 仲村 義平</p> <p>(74) 代理人 100096781 弁理士 堀井 豊</p> <p>(74) 代理人 100098316 弁理士 野田 久登</p>
---	--

最終頁に続く

(54) 【発明の名称】 制御器、制御方法および制御プログラム

(57) 【特許請求の範囲】

【請求項1】

対象とするシステムの時間発展が順方向マルコフ決定過程として記述される際に、前記システムの状態に対する制御則である確率的に表現される方策を前記システムの状態量の観測により方策勾配法によって強化学習する制御器であって、

前記方策に基づいて、前記システムを制御するための制御信号を生成する制御信号生成手段と、

前記システムの前記状態量を観測する状態量検知手段と、

前記状態と前記制御信号とに予め定められた関係で依存する報酬値を獲得する報酬値獲得手段と、

確率的な前記方策を規定するパラメータである方策パラメータにより前記方策が規定される時、各タイムステップにおける前記状態量と前記制御信号とに基づいて、前記システムの状態の分布の定常分布の対数の前記方策パラメータについての偏微分である対数定常分布偏微分を推定することで、前記方策の勾配を推定する方策勾配推定手段と、

前記報酬値と前記方策勾配推定手段による推定結果とに基づいて、前記対数定常分布偏微分を用いて推定した平均報酬偏微分の前記方策パラメータを微小変化させることで、前記方策を更新する方策更新手段とを備える、制御器。

【請求項2】

前記方策勾配推定手段は、前記状態の分布の前記状態についての和が一定であるとの条件により導かれる制約条件であって、前記対数定常分布偏微分の順方向マルコフ連鎖につ

いての期待値が0であるという制約条件の下で、逆方向マルコフ連鎖に対するTD学習により、前記対数定常分布偏微分を推定する、請求項1記載の制御器。

【請求項3】

前記TD学習においては、 i)前記逆方向マルコフ連鎖における前記方策の対数の偏微分の1ステップ前の観測値と1ステップ前の前記対数定常分布偏微分の和と、現在の状態の前記対数定常分布偏微分との差を α とするとき、前記 α の2乗の前記順方向マルコフ連鎖についての期待値と、 i)前記対数定常分布偏微分の前記順方向マルコフ連鎖についての期待値の2乗との和を最小化することにより、前記対数定常分布偏微分を推定する、請求項2記載の制御器。

【請求項4】

対象とするシステムの時間発展が順方向マルコフ決定過程として記述される際に、前記システムの状態に対する制御則である確率的に表現される方策を前記システムの状態量の観測により方策勾配法によって強化学習する制御方法であって、

前記方策に基づいて、前記システムを制御するための制御信号を生成する制御信号生成ステップと、

前記システムの前記状態量を観測する状態量検知ステップと、

前記状態と前記制御信号とに予め定められた関係で依存する報酬値を獲得する報酬値獲得ステップと、

確率的な前記方策を規定するパラメータである方策パラメータにより前記方策が規定される時、各タイムステップにおける前記状態量と前記制御信号とに基づいて、前記システムの状態の分布の定常分布の対数の前記方策パラメータについての偏微分である対数定常分布偏微分を推定することで、前記方策の勾配を推定する方策勾配推定ステップと、

前記報酬値と前記方策勾配推定ステップによる推定結果とに基づいて、前記対数定常分布偏微分に基づき表現される平均報酬の勾配の方向に前記方策パラメータを更新することで、前記方策を更新する方策更新ステップとを備える、制御方法。

【請求項5】

前記方策勾配推定ステップは、前記状態の分布の前記状態についての和が一定であるという条件により導かれる制約条件であって、前記対数定常分布偏微分の順方向マルコフ連鎖についての期待値が0であるという制約条件の下で、逆方向マルコフ連鎖に対するTD学習により、前記対数定常分布偏微分を推定するステップを含む、請求項4記載の制御方法

【請求項6】

前記TD学習においては、 i)前記逆方向マルコフ連鎖における前記方策の対数の偏微分の1ステップ前の観測値と1ステップ前の前記対数定常分布偏微分の和と、現在の状態の前記対数定常分布偏微分との差を α とするとき、前記 α の2乗の前記順方向マルコフ連鎖についての期待値と、 i)前記対数定常分布偏微分の前記順方向マルコフ連鎖についての期待値の2乗との和を最小化することにより、前記対数定常分布偏微分を推定する、請求項5記載の制御方法。

【請求項7】

対象とするシステムの時間発展が順方向マルコフ決定過程として記述される際に、前記システムの状態に対する制御則である確率的に表現される方策を前記システムの状態量の観測により方策勾配法によって強化学習する制御方法をコンピュータに実行させるためのプログラムであって、

前記方策に基づいて、前記システムを制御するための制御信号を生成する制御信号生成ステップと、

前記システムの前記状態量を観測する状態量検知ステップと、

前記状態と前記制御信号とに予め定められた関係で依存する報酬値を獲得する報酬値獲得ステップと、

前記確率的な表現を規定するパラメータである方策パラメータにより前記方策が規定される時、各タイムステップにおける前記状態量と前記制御信号とに基づいて、前記シス

10

20

30

40

50

テムの状態の分布の定常分布の対数の前記方策パラメータについての偏微分である対数定常分布偏微分を推定することで、前記方策の勾配を推定する方策勾配推定ステップと、

前記報酬値と前記方策勾配推定ステップによる推定結果とに基づいて、前記対数定常分布偏微分に基づき表現される平均報酬の勾配の方向に前記方策パラメータを更新することで、前記方策を更新する方策更新ステップとを含む、制御方法をコンピュータに実行させるための制御プログラム。

【請求項 8】

前記方策勾配推定ステップは、前記状態の分布の前記状態についての和が一定であるとの条件により導かれる制約条件であって、前記対数定常分布偏微分の順方向マルコフ連鎖についての期待値が 0 であるという制約条件の下で、逆方向マルコフ連鎖に対する TD 学習により、前記対数定常分布偏微分を推定するステップを含む、請求項 7 記載の制御プログラム。

10

【請求項 9】

前記 TD 学習においては、i) 前記逆方向マルコフ連鎖における前記方策の対数の偏微分の 1 ステップ前の観測値と 1 ステップ前の前記対数定常分布偏微分の和と、現在の状態の前記対数定常分布偏微分との差を とするとき、前記 の 2 乗の前記順方向マルコフ連鎖についての期待値と、i i) 前記対数定常分布偏微分の前記順方向マルコフ連鎖についての期待値の 2 乗との和を最小化することにより、前記対数定常分布偏微分を推定する、請求項 8 記載の制御プログラム。

【発明の詳細な説明】

20

【技術分野】

【0001】

本発明は、方策勾配法により制御対象を制御する制御器、制御方法および制御プログラムの構成に関する。

【背景技術】

【0002】

「マルコフ決定過程」として定式化される制御問題は、ロボット、プラント、移動機械（電車、自動車）などの自律的制御問題として、幅広い応用を持つ重要な技術である。

【0003】

マルコフ決定過程に対する最適制御に関する従来技術として、いわゆる「強化学習」がある。

30

【0004】

「強化学習」とは、エージェントが環境と相互作用を通じて試行錯誤し、得られる累積報酬量を最大化するような「方策」と呼ばれる行動則、すなわち、制御問題に用いる場合には、「制御規則」を学習する理論的な枠組みである。この学習法は、環境やエージェント自身に関する先験的な知識をほとんど必要としない点で様々な分野から注目を集めている。

【0005】

強化学習は大まかに 2 つに分類できる。価値関数を用いて間接的に方策を表現し、価値関数を更新することで方策も更新される「価値関数更新法」と、方策を明示的にもち目的関数の勾配に従って方策を更新する「直接方策更新法（方策勾配法）」である。

40

【0006】

方策勾配法は、行動のランダム性を制御するパラメータも方策パラメータに含めることで確率的方策の獲得が可能であり、また連続系への適用性も高いため、特に注目を集めている。しかし一般に実タスクへ適用すると、適切な行動則を獲得するまでの時間が非現実となることがある。そこで、複数学習器の同時利用、モデルの利用、教示信号の利用等の補助機構を入れて学習時間を短縮させる研究が活発に行われ、成果も著しい。

【0007】

ここで、方策勾配強化学習法（PGRL）は、方策パラメータについての平均報酬の偏微分を用いることにより、方策パラメータを改善して平均報酬を最大化するための強化学

50

習 (RL: Reinforcement Learning) の一般的なアルゴリズムである。ここで、平均報酬の偏微分は、方策勾配 (PG: Policy Gradient) と呼ばれる。従来の方策勾配強化学習法アルゴリズム (PGアルゴリズム) [非特許文献1, 非特許文献2] は、状態の定常分布の偏微分の計算が困難であったため、方策パラメータの変化によりもたらされる定常分布の変化に依存する方策勾配の項を無視している。これらは、定常分布の対数の偏微分 - LSDG (Log Stationary Distribution Gradients) - と呼ばれる。このような省略による偏りは、いわゆる割引率 γ を1に近づければ減少するが、一方で推定された偏微分の分散は多くなってしまふ。このようなトレードオフは、現実には、適切な γ を見いだすことを困難にしていた。マルコフ連鎖の時間を混合する (mixing) ことは適切な γ を決定するための尺度である [非特許文献1, 非特許文献3]。しかしながら、時間の混合は、方策に依存し、したがって、一般には、学習が完了する前に、予測することは困難である。

10

【0008】

また、平均報酬PGアルゴリズム [非特許文献4, 非特許文献5] は、ポアソン方程式の解としての微分コスト関数を導入することにより、割引率を使用しないものである。しかしながら、割引率が1に近い通常のPGと平均報酬PGとのパフォーマンスには、理論的には、大きな相違が無いことが示唆されている [非特許文献6]。

【0009】

したがって、上述したような推定される方策パラメータの偏りを減少させると分散が増大するというトレードオフを解決した学習方法が必要である。

20

【0010】

なお、以下、本文中で引用することとなる方策勾配学習法に関連した先行技術文献を以下に挙げる。

【非特許文献1】Baxter, J. and P. Bartlett (2001) "Infinite-Horizon Policy-Gradient Estimation," *Journal of Artificial Intelligence Research*, Vol. 15, pp. 319-350.

【非特許文献2】H. Kimura and S. Kobayashi. An analysis of actor/critic algorithms using eligibility traces: Reinforcement learning with imperfect value function. In *International Conference on Machine Learning*, 1998.

【非特許文献3】S. Kakade. Optimizing average reward using discounted rewards. In *Annual Conference on Computational Learning Theory*, volume 14. MIT Press, 2001.

30

【非特許文献4】J. N. Tsitsiklis and B. Van Roy. Average cost temporal-difference learning. *Automatica*, 35(11):1799-1808, 1999.

【非特許文献5】V. S. Konda and J. N. Tsitsiklis. On actor-critic algorithms. *SIAM Journal on Control and Optimization*, 42(4):1143-1166, 2001.

【非特許文献6】J. N. Tsitsiklis and B. Van Roy. On average versus discounted reward temporal-difference learning. *Machine Learning*, 49(2):179-191, 2002.

【非特許文献7】P. W. Glynn. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33(10):75-84, 1991.

【非特許文献8】R. Y. Rubinstein. How to optimize discrete-event system from a single sample path by the score function method. *Annals of Operations Research*, 27(1):175-212, 1991.

40

【非特許文献9】A. Y. Ng, R. Parr, and D. Koller. Policy search via density estimation. In *Advances in Neural Information Processing Systems*. MIT Press, 2000.

【非特許文献10】D. P. Bertsekas. *Dynamic Programming and Optimal Control*, Volumes 1 and 2. Athena Scientific, 1995.

【非特許文献11】R. S. Sutton and A. G. Barto. *Reinforcement Learning*. MIT Press, 1998.

【非特許文献12】R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*, 3:9-44, 1988.

【非特許文献13】S. J. Bradtke and A. G. Barto. Linear least-squares algorithms

50

for temporal difference learning. Machine Learning, 22(1-3):33-57, 1996.

【非特許文献14】J. A. Boyan. Least-squares temporal difference learning. Machine Learning, 49(2-3):233-246, 2002.

【非特許文献15】R. B. Schinazi. Classical and Spatial Stochastic Processes. Birkhauser, 1999.

【非特許文献16】J. Peng and R. J. Williams. Incremental multi-step Q-learning. Machine Learning, 22:283-290, 1996.

【非特許文献17】P. Young. Recursive Estimation and Time-series Analysis. Springer-Verlag, 1984.

【非特許文献18】D. P. Bertsekas and J. N. Tsitsiklis. Neuro-Dynamic Programming. Athena Scientific, 1996. 10

【発明の開示】

【発明が解決しようとする課題】

【0011】

すでに、(定常)状態分布の勾配を評価するための2つの方法があるが、これらは、本願発明の方法とは異なったものであり、以下のような問題がある。最初のものは、「尤度比勾配法」あるいは「スコア関数法」と呼ばれるものであり[非特許文献7, 非特許文献8]、再生プロセスにしか適用できない問題がある[非特許文献1]。非特許文献9に開示された他の方法は、状態分布の直接的な評価ではなく、密度伝搬を伴う状態分布の評価を介して実行されるものである。したがって、これらの方法は、エージェントがどの状態にあるかの知識を必要とするのに対し、後に説明する本発明の方法では、ノイズを含む状態の特徴ベクトルを観測するのみでよい。 20

【0012】

したがって、本発明は、上記のような問題点を解決するためになされたものであって、その目的は、方策パラメータの偏りを減少させると分散が増大するというトレードオフを解決した方策勾配学習方法を用いた制御器、制御方法または制御プログラムを提供することである。

【0013】

本発明の他の目的は、汎用性のある方策勾配学習方法を用いた制御器、制御方法または制御プログラムを提供することである。 30

【課題を解決するための手段】

【0014】

このような目的を達成するために、本発明の制御器は、対象とするシステムの時間発展が順方向マルコフ決定過程として記述される際に、システムの状態に対する制御則である確率的に表現される方策をシステムの状態量の観測により方策勾配法によって強化学習する制御器であって、方策に基づいて、システムを制御するための制御信号を生成する制御信号生成手段と、システムの状態量を観測する状態量検知手段と、状態と制御信号とに予め定められた関係で依存する報酬値を獲得する報酬値獲得手段と、確率的な方策を規定するパラメータである方策パラメータにより方策が規定されるとき、各タイムステップにおける状態量と制御信号とに基づいて、システムの状態の分布の定常分布の対数の方策パラメータについての偏微分である対数定常分布偏微分を推定することで、方策の勾配を推定する方策勾配推定手段と、報酬値と方策勾配推定手段による推定結果とに基づいて、対数定常分布偏微分を用いて推定した平均報酬偏微分の方向に方策パラメータを微小変化させることで、方策を更新する方策更新手段とを備える。 40

【0015】

好ましくは、方策勾配推定手段は、状態の分布の状態についての和が一定であるとの条件により導かれる制約条件であって、対数定常分布偏微分の順方向マルコフ連鎖についての期待値が0であるという制約条件の下で、逆方向マルコフ連鎖に対するTD学習により、対数定常分布偏微分を推定する。

好ましくは、TD学習においては、i) 逆方向マルコフ連鎖における方策の対数の偏微 50

分の1ステップ前の観測値と1ステップ前の対数定常分布偏微分の和と、現在の状態の対数定常分布偏微分との差をとするとき、 θ の2乗の順方向マルコフ連鎖についての期待値と、 i の対数定常分布偏微分の順方向マルコフ連鎖についての期待値の2乗との和を最小化することにより、対数定常分布偏微分を推定する。

【0016】

この発明の他の局面に従うと、対象とするシステムの時間発展が順方向マルコフ決定過程として記述される際に、システムの状態に対する制御則である確率的に表現される方策をシステムの状態量の観測により方策勾配法によって強化学習する制御方法であって、方策に基づいて、システムを制御するための制御信号を生成する制御信号生成ステップと、システムの状態量を観測する状態量検知ステップと、状態と制御信号とに予め定められた関係で依存する報酬値を獲得する報酬値獲得ステップと、確率的な方策を規定するパラメータである方策パラメータにより方策が規定されるとき、各タイムステップにおける状態量と制御信号とに基づいて、システムの状態の分布の定常分布の対数の方策パラメータについての偏微分である対数定常分布偏微分を推定することで、方策の勾配を推定する方策勾配推定ステップと、報酬値と方策勾配推定手段による推定結果とに基づいて、対数定常分布偏微分に基づき表現される平均報酬の勾配の方向に方策パラメータを更新することで、方策を更新する方策更新ステップとを備える。

【0017】

この発明のさらに他の局面に従うと、対象とするシステムの時間発展が順方向マルコフ決定過程として記述される際に、システムの状態に対する制御則である確率的に表現される方策をシステムの状態量の観測により方策勾配法によって強化学習する制御方法をコンピュータに実行させるためのプログラムであって、方策に基づいて、システムを制御するための制御信号を生成する制御信号生成ステップと、システムの状態量を観測する状態量検知ステップと、状態と制御信号とに予め定められた関係で依存する報酬値を獲得する報酬値獲得ステップと、確率的な表現を規定するパラメータである方策パラメータにより方策が規定されるとき、各タイムステップにおける状態量と制御信号とに基づいて、システムの状態の分布の定常分布の対数の方策パラメータについての偏微分である対数定常分布偏微分を推定することで、方策の勾配を推定する方策勾配推定ステップと、報酬値と方策勾配推定ステップによる推定結果とに基づいて、対数定常分布偏微分に基づき表現される平均報酬の勾配の方向に方策パラメータを更新することで、方策を更新する方策更新ステップとを含む、制御方法をコンピュータに実行させる。

【発明を実施するための最良の形態】

【0018】

以下の説明の構成の概要を説明すると、(1. 本発明の概要)において、本発明の全体的な構成を説明し、(2. 前提：方策勾配強化学習)では、従来のPG法をレビューし、本発明のLSDGを評価する目的を概説する。(3. 定常分布の対数の偏微分の推定)では、逆方向マルコフ連鎖の方法に基づく最小二乗TD法によるLSDG()アルゴリズムを説明する。(4. LSDG推定による方策の更新)では、LSDG()-PGアルゴリズムを説明する。これは、LSDG()を利用し、割引率 γ を用いない方法である。(5. 数値計算の結果)では、提案する方法のパフォーマンスを確認するため、簡単なマルコフ決定過程(MDP)における計算結果が示される。

【0019】

(1. 本発明の概要)

後に説明するように、本発明では、逆方向マルコフ連鎖の方法とTD(temporal difference)学習アルゴリズムにより、定常分布の勾配としてのLSDGが導出される、新たな方策勾配法の枠組みを提案する。この枠組みにおいては、平均報酬の偏微分は、割引率に依存せず、 $\gamma = 0$ とおくことにより、価値関数を学習する必要がなくなる。

【0020】

以下、図面を参照して本発明の実施の形態について説明する。

以下の説明で明らかとなるとおり、本発明は、ロボット、プラント、移動機械(電車、

10

20

30

40

50

自動車)などの制御問題として、幅広い応用を持つ。

【0021】

ただし、以下では、本発明の具体的な適用例として、特に簡単なロボットの自動制御問題を対象とするものとして説明を行う。また、数値計算の結果は、さらに簡単なモデルに対する比較を示している。しかしながら、本発明は、このような応用に限定されるものではなく、より一般的に、対象システムの時間発展が複雑な場合の対象システムの制御に適用することができる。そのようなものの例としては、巨大プラント(溶鉱炉、原子力プラント)、マルチリンクロボット(ヒューマノイドロボット)、ノンホロノームシステム(宇宙ステーション)、地下鉄ホームでの人の流れなどがある。これらは、いずれも古典的制御法での制御が困難であり、かつ重要な制御対象である。

10

【0022】

(1. 本発明のシステム構成)

図1は、本発明の制御方法および制御プログラムが適用される制御器を用いたシステム1000の一例を示す概念図である。

【0023】

図1を参照して、システム1000は、制御対象となる被制御装置200と、この被制御装置200に対して制御信号を与えるためのコンピュータ100とを備える。

【0024】

図1を参照してこのコンピュータ100は、CD-ROM(Compact Disc Read-Only Memory)上の情報を読み込むためのCD-ROMドライブ108およびフレキシブルディスク(Flexible Disk、以下FD)116に情報を読み書きするためのFDドライブ106を備えたコンピュータ本体102と、コンピュータ本体102に接続された表示装置としてのディスプレイ104と、同じくコンピュータ本体102に接続された入力装置としてのキーボード110およびマウス112とを含む。

20

【0025】

図2は、このコンピュータ100の構成をブロック図形式で示す図である。

図2に示されるように、このコンピュータ100を構成するコンピュータ本体102は、CD-ROMドライブ108およびFDドライブ106に加えて、それぞれバスBSに接続されたCPU(Central Processing Unit)120と、ROM(Read Only Memory)およびRAM(Random Access Memory)を含むメモリ122と、直接アクセスメモリ装置、たとえば、ハードディスク124と、被制御装置200とデータの授受を行うための通信インタフェース128とを含んでいる。CD-ROMドライブ108にはCD-ROM118が装着される。FDドライブ106にはFD116が装着される。

30

【0026】

被制御装置200からは、コンピュータ100に対して被制御装置200の状態を示すパラメータ(状態量)の情報、たとえば、被制御装置200の可動部分の位置、速度、加速度、角度、角速度等の情報が与えられる。一方、コンピュータ100からは、被制御装置200に対して、これら状態量を制御するための制御情報が制御信号として与えられる。

【0027】

なお、CD-ROM118は、コンピュータ本体に対してインストールされるプログラム等の情報を記録可能な媒体であれば、他の媒体、たとえば、DVD-ROM(Digital Versatile Disc)やメモリカードなどでもよく、その場合は、コンピュータ本体102には、これらの媒体を読み取ることが可能なドライブ装置が設けられる。

40

【0028】

本発明の制御器の主要部は、コンピュータハードウェアと、CPU120により実行されるソフトウェアとにより構成される。一般的にこうしたソフトウェアはCD-ROM118、FD116等の記憶媒体に格納されて流通し、CD-ROMドライブ108またはFDドライブ106等により記憶媒体から読取られてハードディスク124に一旦格納される。または、当該装置がネットワークに接続されている場合には、ネットワーク上のサ

50

ーバから一旦ハードディスク 1 2 4 にコピーされる。そしてさらにハードディスク 1 2 4 からメモリ 1 2 2 中の R A M に読出されて C P U 1 2 0 により実行される。なお、ネットワーク接続されている場合には、ハードディスク 1 2 4 に格納することなく R A M に直接ロードして実行するようにしてもよい。

【 0 0 2 9 】

図 1 および図 2 に示したコンピュータのハードウェア自体およびその動作原理は一般的なものである。したがって、本発明の最も本質的な部分は、F D 1 1 6、C D - R O M 1 1 8、ハードディスク 1 2 4 等の記憶媒体に記憶されたソフトウェアである。

【 0 0 3 0 】

なお、一般的傾向として、コンピュータのオペレーティングシステムの一部として様々なプログラムモジュールを用意しておき、アプリケーションプログラムはこれらモジュールを所定の配列で必要な時に呼び出して処理を進める方式が一般的である。そうした場合、当該制御器を実現するためのソフトウェア自体にはそうしたモジュールは含まれず、当該コンピュータでオペレーティングシステムと協働してはじめて制御器が実現することになる。しかし、一般的なプラットフォームを使用する限り、そうしたモジュールを含ませたソフトウェアを流通させる必要はなく、それらモジュールを含まないソフトウェア自体およびそれらソフトウェアを記録した記録媒体（およびそれらソフトウェアがネットワーク上を流通する場合のデータ信号）が実施の形態を構成すると考えることができる。

【 0 0 3 1 】

[制御方法の一般的説明]

以下、本発明の構成について、その理論的な構成をまず説明する。

【 0 0 3 2 】

(制御器の構成)

(2 . 前提 : 方策勾配強化学習)

以下では、マルコフ決定過程 (M D P) について考えることにし、制御対象 (制御器の環境) は状態遷移確率 (時刻 t において、状態 x_t であるときに行動 (制御) u_t を実行することで状態が x_{t+1} となる確率) と報酬関数 $r_{t+1} = r (x_t, u_t)$ (なお、 $r_{t+1} = r (x_{t+1}, x_t, u_t)$ の場合にも同様に議論できる) によって特徴づけられるものとする。なお、この報酬関数については、制御対象の制御目標に応じて予め定められているものとする。状態入力 $x \in X$ から行動出力 $u \in U$ への写像を方策と呼び、以下で説明するように確率的に表現する。方策は、パラメータ π で、規定される。

【 0 0 3 3 】

ここでは、従来の P G R L アルゴリズムをレビューし、本発明の新しいアルゴリズムの主たるアイデアを提示する。有限な状態 $X = \{x\}$ の組と行動 $U = \{u\}$ とを有する離散時間 M D P は、以下の状態遷移確率 p と報酬関数 r_{+1} によって規定される。

【 0 0 3 4 】

【 数 1 】

$$p(x_{+1} | x, u)$$

$$r_{+1} = r(x_{+1}, x, u) \in [R_{\min}, R_{\max}]$$

【 0 0 3 5 】

ここで、記載の簡単のために、 x_{+1} は、状態 x において、行動 u により与えられる次の状態であり、 r_{+1} は、 x_{+1} において観測された即時報酬である [非特許文献 1 0 , 非特許文献 1 1]。 x_{+k} および u_{+k} は、それぞれ、状態 x から k 時間ステップ先の状態および行動であり、添え字が $-k$ となっていれば、その反対である。(M D P における) 決定は、 R^d によりパラメータ化された以下の確率の方策 π にしたがってなされる。

【 0 0 3 6 】

10

20

30

40

【数 2】

$$\pi_{\theta}(x, u) \equiv p(\mathbf{u} | \mathbf{x}; \theta)$$

【0037】

以下の方策 π は、全ての $x \in X$ および $u \in U$ に対して、 π について微分可能であると仮定する。

【0038】

【数 3】

$$\pi_{\theta}(x, u)$$

10

【0039】

さらに、以下のような仮定をおく。

(仮定 1)

以下のような状態遷移確率 p と確率の方策 π とを有するマルコフ連鎖 $M(\pi)$ は、全ての方策パラメータについてエルゴード的 (既約で非周期的) である。

【0040】

【数 4】

$$p(x_{+1} | x, u)$$

20

$$\pi_{\theta}(x, u) \equiv p(\mathbf{u} | \mathbf{x}; \theta)$$

【0041】

したがって、以下のただ 1 つの定常分布が存在する。

【0042】

【数 5】

$$d^{\pi}(\mathbf{x}) = \lim_{k \rightarrow \infty} p(x_{+k} = x | \pi_{\theta}) > 0$$

【0043】

この定常分布は、初期状態には独立であって、以下の式を満たす。

30

【0044】

【数 6】

$$d^{\pi}(\mathbf{x}) = \sum_{x_{-1}, u_{-1}} p_{M(\theta)}(x, u_{-1} | x_{-1}) d^{\pi}(x_{-1}) \quad (1)$$

【0045】

ここで、以下の式が成り立つ。

【0046】

【数 7】

$$p_{M(\theta)}(x, u_{-1} | x_{-1}) \equiv \pi_{\theta}(x_{-1}, u_{-1}) p(x | x_{-1}, u_{-1})$$

40

【0047】

PRGL の目的は、以下の「平均報酬」と呼ばれる即時報酬の平均を最大化する方策パラメータ θ^* を見いだすことである。

【0048】

【数 8】

$$R(\theta) \equiv \lim_{K \rightarrow \infty} \frac{1}{K} E_{M(\theta)} \left\{ \sum_{k=1}^K r_{+k} \mid x \right\}$$

ここで、 $E_{M(\theta)}$ はマルコフ連鎖 $M(\theta)$ についての期待値である。

【0049】

仮定 1 の下では、平均報酬は、初期状態 x には独立で、以下の式に等しいことが示される [非特許文献 10] : 10

【0050】

【数 9】

$$R(\theta) = \sum_{x_{+1}, x, u} d^\pi(x) \pi_\theta(x, u) p(x_{+1} \mid x, u) r(x_{+1}, x, u) \quad (2)$$

【0051】

方策勾配 RL アルゴリズムは、方策パラメータ θ を、以下の式に示される $R(\theta)$ についての平均報酬 $R(\theta)$ の勾配の方向に更新する。

【0052】

【数 10】

$$\nabla_\theta R(\theta) \equiv \left[\frac{\partial}{\partial \theta_1} R(\theta), \dots, \frac{\partial}{\partial \theta_d} R(\theta) \right]^T$$

【0053】

以下、単に方策勾配 (PG) と、しばしば呼ばれる。この方策勾配は、以下のように与えられる。

【0054】

【数 11】

$$\nabla_\theta R(\theta) = \sum_{x_{+1}, x, u} d^\pi(x) \pi_\theta(x, u) (\nabla_\theta \ln \pi_\theta(x, u) + \nabla_\theta \ln d^\pi(x)) p(x_{+1} \mid x, u) r(x_{+1}, x, u) \quad (3)$$

【0055】

以下の式に示される定常分布の対数の偏微分の導出は簡単ではない。

【0056】

【数 12】

$$\nabla_\theta \ln d^\pi(x)$$

【0057】

そこで、従来の PG アルゴリズム [非特許文献 1, 非特許文献 2] は、PG のもう一つの表現を利用している。

【0058】

10

20

30

40

【数 1 3】

$$\nabla_{\theta} R(\theta) = \sum_{x,u} d^{\pi}(x) \pi_{\theta}(x,u) \nabla_{\theta} \ln \pi_{\theta}(x,u) Q_{\gamma}^{\pi}(x,u) + (1-\gamma) \sum_x d^{\pi}(x) \nabla_{\theta} \ln d^{\pi}(x) V_{\gamma}^{\pi}(x) \quad (4)$$

【0059】

ここで、割引率 $\gamma \in [0, 1)$ で、それぞれ、行動価値関数 Q と状態価値関数 V とは以下のように表される [非特許文献 11]。

【0060】

【数 1 4】

$$Q_{\gamma}^{\pi}(x,u) \equiv \lim_{K \rightarrow \infty} E_{M(\theta)} \left\{ \sum_{k=1}^K \gamma^{k-1} r_{x+u}^k \mid x,u \right\}$$

$$V_{\gamma}^{\pi}(x,u) \equiv \lim_{K \rightarrow \infty} E_{M(\theta)} \left\{ \sum_{k=1}^K \gamma^{k-1} r_{x+u}^k \mid x \right\}$$

【0061】

式(4)の第2項の寄与は、 γ が1に近づくにつれて小さくなるので [非特許文献 1]、従来のアルゴリズム [非特許文献 1, 非特許文献 2] は、 $\gamma \sim 1$ とすることで、第1項のみから PG を近似している。このような省略による偏りは、割引率 γ を1に近づければ小さくなるが、推定の分散は多くなってしまう。

【0062】

ここで、本発明では、もう1つのアプローチを提案する。そこでは、以下の式の定常分布の対数の偏微分 (LSDG) を推定し、PG を導出するために式(3)を用いる。

【0063】

【数 1 5】

$$\nabla_{\theta} \ln d^{\pi}(x)$$

【0064】

著しい特徴は、価値関数を学習する必要がなく、したがって、そのアルゴリズムは、割引率 γ の選択において、偏りと分散のトレードオフと関係がないことである。

(3. 定常分布の対数の偏微分の推定)

以下では、最小二乗法に基づく LSDG 推定アルゴリズム、LSDG () を提案する。この目的のために、エルゴード的なマルコフ連鎖 M () の逆過程を定式化し、LSDG は、TD法 [非特許文献 12, 非特許文献 13, 非特許文献 14] の枠組みで推定できることを示す。

(3.1 順方向および逆方向マルコフ連鎖の性質)

ベイズの理論を用いれば、現在の状態から過去の状態および行動の対への逆方向の確率は、以下の式で表される。

【0065】

【数 1 6】

$$q(x_{-1}, u_{-1} \mid x) = \frac{p(x \mid x_{-1}, u_{-1}) p(x_{-1}, u_{-1})}{\sum_{x_{-1}, u_{-1}} p(x \mid x_{-1}, u_{-1}) p(x_{-1}, u_{-1})}$$

【0066】

以下の事後確率 q は事前分布 p に依存する。

【0067】

【数 1 7】

$$q(x_{-1}, u_{-1} \mid x)$$

$$p(x_{-1}, u_{-1})$$

10

20

30

40

50

【 0 0 6 8 】

以下のように、事前分布 p が定常分布 d と方策 に従うとき、事後分布 q は、定常逆方向確率と呼ばれ、下付添え字 B () が加えられる。

【 0 0 6 9 】

【 数 1 8 】

$$p(x_{-1}, u_{-1}) = \pi_{\theta}(x_{-1}, u_{-1}) d^{\pi}(x_{-1})$$

ここで、 $q_{B(\theta)}(x_{-1}, u_{-1} | x)$ は、以下のように表される：

10

$$\begin{aligned} q_{B(\theta)}(x_{-1}, u_{-1} | x) &= \frac{p(x | x_{-1}, u_{-1}) \pi_{\theta}(x_{-1}, u_{-1}) d^{\pi}(x_{-1})}{d^{\pi}(x)} \\ &= \frac{p_{M(\theta)}(x, u_{-1} | x_{-1}) d^{\pi}(x_{-1})}{d^{\pi}(x)} \end{aligned} \quad (5)$$

もしも、マルコフ連鎖が $q_{B(\theta)}(x_{-1}, u_{-1} | x)$ に従うならば、それを $p_{M(\theta)}(x, u_{-1} | x_{-1})$ に従う順方向マルコフ連鎖 $M(\theta)$ に関連する逆方向マルコフ連鎖 $B(\theta)$ と呼ぶ。

20

【 0 0 7 0 】

マルコフ連鎖 M () と B () の両方 は、以下の 2 つの定理において記述されるように密接に関連している。

(定理 1)

【 0 0 7 1 】

【 数 1 9 】

30

遷移確率 $p_{M(\theta)}(x | x_{-1}) \equiv \sum_{u_{-1}} p_{M(\theta)}(x, u_{-1} | x_{-1})$ により特徴づけられるマルコフ連鎖 $M(\theta)$ が既約でエルゴード的であるとき、 $p_{M(\theta)}(x | x_{-1})$ に対する逆方向 (定常) 遷移確率 $q_{B(\theta)}(x_{-1} | x) \equiv \sum_{u_{-1}} q_{B(\theta)}(x_{-1}, u_{-1} | x)$ により特徴づけられる逆方向マルコフ連鎖 $B(\theta)$ もまた、エルゴード的であり、 $M(\theta)$ と同一のただ一つの定常分布を有する：

40

$$d_{M(\theta)}(x) = d_{B(\theta)}(x) \quad (6)$$

ここで、 $d_{M(\theta)}(x) \equiv d^{\pi}(x)$ と $d_{B(\theta)}(x)$ は、それぞれ、 $M(\theta)$ と $B(\theta)$ の定常分布である。

【 0 0 7 2 】

(証明)

式 (5) の両辺に以下の定常分布をかける。

【 0 0 7 3 】

50

【数 2 0】

$$d^\pi(x)$$

【0 0 7 4】

すると、全ての可能な行動 $u_{-1} \in U$ について総和をとると、以下の式が得られる：

【0 0 7 5】

【数 2 1】

$$q_{B(\theta)}(x_{-1}|x)d^\pi(x) = P_{M(\theta)}(x|x_{-1})d^\pi(x_{-1}) \quad (7)$$

【0 0 7 6】

そして、式 (7) の両辺を可能な状態 $x \in X$ について総和をとると、以下の式が成り立つ。

【0 0 7 7】

【数 2 2】

$$\sum_x q_{B(\theta)}(x_{-1}|x)d^\pi(x) = d^\pi(x_{-1})$$

【0 0 7 8】

このことは、以下の 2 点を成立させる。(i) $B(\theta)$ は、 $M(\theta)$ と同一の定常分布を有すること、(ii) $B(\theta)$ は $M(\theta)$ と同じ既約な性質を有すること。

【0 0 7 9】

【数 2 3】

式 (7) は、遷移確率 $P_{M(\theta)}(x|x_{-1})$ または $p_{B(\theta)}(x_{-1}|x)$ により、行列表現 $P_{M(\theta)}$ また

は $Q_{B(\theta)}$ にそれぞれまとめ、定常分布をベクトル表現 d^π にまとめると：

$$Q_{B(\theta)} = \text{diag}(d^\pi)^{-1} P_{M(\theta)}^T \text{diag}(d^\pi)$$

容易に、任意の自然数 n に対して $(P_{M(\theta)})^n$ の対角成分は、 $(Q_{B(\theta)})^n$ のそれに等しいこ

とがわかる。

【0 0 8 0】

このことは、(iii) $B(\theta)$ が $M(\theta)$ と同じ非周期的な性質をもっていることを示唆する。式 (6) は、(i) - (iii) により直接証明される [非特許文献 15]。

(定理 2)

【0 0 8 1】

10

20

30

【数 2 4】

x_{-K} の分布が $d^\pi(x)$ にしたがうとき、 $k \in [0, K]$ にわたる任意の関数 $f(x_{-k}, u_{-k})$

の和に関する順方向逆方向のマルコフ連鎖の期待値は、相互に等しい：

$$\mathbf{E}_{B(\theta)} \left\{ \sum_{k=0}^K f(x_{-k}, u_{-k}) \mid x \right\} = \mathbf{E}_{M(\theta)} \left\{ \sum_{k=0}^K f(x_{-k}, u_{-k}) \mid x, d^\pi(x_{-K}) \right\}$$

(8)

10

ここで、 $\mathbf{E}_{B(\theta)}$ および $\mathbf{E}_{M(\theta)}$ は、順方向および逆方向マルコフ連鎖 $B(\theta)$ と $M(\theta)$

とのそれぞれの期待値を示し、

$$\mathbf{E} \left\{ \mid d^\pi(x_{-K}) \right\} \equiv \mathbf{E} \left\{ \mid p(x_{-K}) = d^\pi(x_{-K}) \right\}$$

である。式 (8) は、 $K \rightarrow \infty$ の極限においても成り立つ。

【0082】

20

(証明)

マルコフ連鎖の特性と式 (5) を代入することにより、以下の関係が得られる。

【0083】

【数 2 5】

$$\begin{aligned} & q_{B(\theta)}(x_{-1}, u_{-1}, \dots, x_{-K}, u_{-K} \mid x) \\ &= q_{B(\theta)}(x_{-1}, u_{-1} \mid x) \cdots q_{B(\theta)}(x_{-K}, u_{-K} \mid x_{-K+1}) \\ &\propto p_{M(\theta)}(x, u_{-1} \mid x_{-1}) \cdots p_{M(\theta)}(x_{-K+1}, u_{-K} \mid x_{-K}) d^\pi(x_{-K}) \end{aligned}$$

30

【0084】

このことは、有限の K の場合において式 (8) が成立することを証明する。定理 1 から以下の式が導かれるので、式 (8) の K の極限の場合も成立することが、すぐさま証明される。

【0085】

【数 2 6】

$$\begin{aligned} \lim_{K \rightarrow \infty} \mathbf{E}_{B(\theta)} \{ f(x_{-k}, u_{-k}) \mid x \} &= \lim_{K \rightarrow \infty} \mathbf{E}_{M(\theta)} \{ f(x_{-k}, u_{-k}) \mid x, d^\pi(x_{-K}) \} \\ &= \sum_{x,u} \pi_\theta(x,u) d^\pi(x) f(x,u) \end{aligned}$$

40

【0086】

定理 1 および定理 2 は、これらが、定常分布に収束する状態分布の下で、順方向マルコフ連鎖 $M(\quad)$ からのサンプルは、そのまま、逆方向マルコフ連鎖 $B(\quad)$ に関する推定に使用できることになるので、重要である。そして、これらは、後に利用されるものである。

(3.2 逆方向から順方向のマルコフ連鎖の LSDG のための TD (Temporal Difference) 学習法)

LSDG は式 (5) を用いて、以下のように分解される。

【0087】

50

【数 27】

$$\begin{aligned}
\nabla_{\theta} \ln d^{\pi}(x) &= \frac{1}{d^{\pi}(x)} \sum_{x_{-1}, u_{-1}} p(x | x_{-1}, u_{-1}) \pi_{\theta}(x_{-1}, u_{-1}) d^{\pi}(x_{-1}) \\
&\quad \left\{ \nabla_{\theta} \ln \pi_{\theta}(x_{-1}, u_{-1}) + \nabla_{\theta} \ln d^{\pi}(x_{-1}) \right\} \\
&= \sum_{x_{-1}, u_{-1}} q_{B(\theta)}(x_{-1}, u_{-1} | x) \left\{ \nabla_{\theta} \ln \pi_{\theta}(x_{-1}, u_{-1}) + \nabla_{\theta} \ln d^{\pi}(x_{-1}) \right\} \\
&= \mathbf{E}_{B(\theta)} \left\{ \nabla_{\theta} \ln \pi_{\theta}(x_{-1}, u_{-1}) + \nabla_{\theta} \ln d^{\pi}(x_{-1}) \mid x \right\}
\end{aligned} \tag{9}$$

ここで、式(9)において $\nabla_{\theta} \ln d^{\pi}(x)$ および $\nabla_{\theta} \ln d^{\pi}(x_{-1})$ が存在することに注意して、式(9)の反復により、以下の式が得られる。

$$\nabla_{\theta} \ln d^{\pi}(x) = \lim_{K \rightarrow \infty} \mathbf{E}_{B(\theta)} \left\{ \sum_{k=1}^K \nabla_{\theta} \ln \pi_{\theta}(x_{-k}, u_{-k}) + \nabla_{\theta} \ln d^{\pi}(x_{-K}) \mid x \right\} \tag{10}$$

【0088】

式(10)は、状態 x のLSDGは、以下の式で表される方策の対数の偏微分の状態 x から逆方向マルコフ連鎖 $B(\quad)$ の無限区間の集積であることを暗示している。

【0089】

【数28】

$$\nabla_{\theta} \ln \pi_{\theta}(x, u)$$

【0090】

式(9)および(10)から、LSDGは、 $M(\quad)$ よりもむしろ逆方向マルコフ連鎖 $B(\quad)$ についての、以下のような逆方向TDに関するTD学習[非特許文献12]により推定されうる。

【0091】

【数29】

$$\delta(x) \equiv \nabla_{\theta} \ln \pi_{\theta}(x_{-1}, u_{-1}) + \nabla_{\theta} \ln d^{\pi}(x_{-1}) - \nabla_{\theta} \ln d^{\pi}(x)$$

【0092】

ここで、最初の2つの項は、 $B(\quad)$ における方策の対数の偏微分の1ステップ前の実際の観測値と1ステップ前のLSDGであり、現在の状態のLSDGが最後の項である。

【0093】

10

... (9)

20

... (10)

30

40

【数30】

$\delta(x)$ はランダム変数である一方、 $\mathbf{E}_{B(\theta)}\{\delta(x)|x\} = 0$ が成り立つ。LSDG の推定のために、逆方向TDエラーの二乗の平均 $\mathbf{E}_{B(\theta)}\left\{\hat{\delta}(x)^2\right\}$ を最小化することになる。ここで、 $\hat{\delta}(x)$ は、LSDG $\nabla_{\theta} \ln d^{\pi}(x)$ よりもむしろ LSDG 推定値 $\hat{\nabla}_{\theta} \ln d^{\pi}(x)$ によって構成される。 $\hat{\delta}(x)^2$ は、 $\delta(x)^T \delta(x)$ を簡単に示すものである。

10

【0094】

適格度減衰率 $[0, 1]$ と逆方向追跡時間ステップ K N を用いて、式(10)は、以下のように一般化される。

【0095】

【数31】

$$\begin{aligned} & \nabla_{\theta} \ln d^{\pi}(x) \\ &= \mathbf{E}_{B(\theta)} \left\{ \sum_{k=1}^K \lambda^{k-1} \left\{ \nabla_{\theta} \ln \pi_{\theta}(x_{-k}, u_{-k}) + (1-\lambda) \nabla_{\theta} \ln d^{\pi}(x_{-K}) \right\} + \lambda^k \nabla_{\theta} \ln d^{\pi}(x_{-K}) | x \right\} \end{aligned}$$

20

このような変形に従って、逆方向TDは、逆方向TD $(\lambda), \delta_{\lambda, K}(x)$ に変形される。

$$\begin{aligned} \delta_{\lambda, K}(x) \equiv & \sum_{k=1}^K \lambda^{k-1} \left\{ \nabla_{\theta} \ln \pi_{\theta}(x_{-k}, u_{-k}) + (1-\lambda) \nabla_{\theta} \ln d^{\pi}(x_{-K}) \right\} \\ & + \lambda^k \nabla_{\theta} \ln d^{\pi}(x_{-K}) - \nabla_{\theta} \ln d^{\pi}(x) \end{aligned}$$

30

ここで、偏り無しの特性的、 $\mathbf{E}_{B(\theta)}\{\delta_{\lambda, K}(x)|x\} = 0$ は、依然、保持される。 $\lambda =$

1かつ $K \rightarrow \infty$ の極限において、 $\mathbf{E}_{B(\theta)}\left\{\hat{\delta}_{\lambda, K}(x)^2\right\}$ を最小化することは、ウィドロ

40

ーホッフ (Widrow-Hoff) 教師付き学習とみなされる。

【0096】

上記のような設定でなくとも、大きな λ や K を用いたならば、このような最小化は従来の価値関数に対するTD()学習の場合のように、非マルコフ効果に対しては、より敏感ではない[非特許文献16]。

【0097】

【数 3 2】

LSDGの推定として $\mathbf{E}_{B(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x)^2 \right\}$ を最小化するために、逆方向マルコフ

連鎖 $B(\theta)$ から取り出された多くのサンプルを集める必要がある。幸いなことに、定理 1 および定理 2 を利用すると、以下のように交換可能な特性を用いることができる：

$$\mathbf{E}_{B(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x)^2 \right\} = \sum_x d_{B(\theta)}(x) \mathbf{E}_{B(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x)^2 \mid x \right\} \quad 10$$

$$= \sum_x d^\pi(x) \mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x)^2 \mid x, d^\pi(x_{-K}) \right\}$$

$$= \mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x)^2 \mid d^\pi(x_{-K}) \right\}$$

… (11)

20

すなわち、 $x_{-K} \sim d^\pi(x)$ であれば、 $\mathbf{E}_{B(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x)^2 \right\}$ を最小化するために、

実際のサンプルを再利用できる。現実の問題では、しかしながら、初期状態は、定常状態 $d^\pi(x)$ から取り出されることはまれである。

【0098】

理論的な仮定と現実への適用との間のギャップを埋めるために、以下の2つのうちのいずれかの戦略をとる必要がある。(i) ~ 1 ならば、 K は、あまり大きな整数に設定しない、(ii) $K \sim t$ ならば、 ~ 1 に設定しない、ここで、 t は、現実の順方向マルコフ連鎖の現在のタイムステップである。

30

(3.3 LSDG推定アルゴリズム：制限付き逆方向TD()の最小二乗法)

【0099】

【数 3 3】

3.2では、LSDGの推定がマルコフ連鎖 $M(\theta)$ における $\hat{\delta}_{\lambda, K}(x)^2$ の平均二

乗 $\mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x)^2 \mid d^\pi(x_{-K}) \right\}$ を最小化することにより実行できるという

40

理論を紹介した。しかしながら、LSDGには、また、 $\sum_x d^\pi(x) = 1$ との条件から導かれる以下のような制約条件がある。

$$\mathbf{E}_{M(\theta)} \left\{ \nabla_\theta \ln d^\pi(x) \right\} = \sum_x d^\pi(x) \nabla_\theta \ln d^\pi(x) = \nabla_\theta \sum_x d^\pi(x) = 0$$

… (12)

50

【 0 1 0 0 】

この 3.3 では、最小二乗法に基づく [非特許文献 17 , 非特許文献 13 , 非特許文献 14]、LSDG 推定アルゴリズム、LSDG () を提案する。これは、同時に、平均二乗を減少させるとともに、制約条件を満足することを達成しようとするものである。

【 0 1 0 1 】

【 数 3 4 】

LSDG 推定 $\hat{\nabla}_{\theta} \ln d^{\pi}(x)$ が、線形ベクトル関数近似器 $f(x; \Omega) \equiv \Omega \phi(x)$ で表されると考える。ここで、 $\phi(x) \in R^e$ は基底関数であり、
 $\Omega \equiv [\omega_1, \dots, \omega_d]^T \in R^{d \times e}$ は、調整可能なパラメータ行列であり、最適パラメータ Ω^* は、 $\nabla_{\theta} \ln d^{\pi}(x) = \Omega^* \phi(x)$ を満たす。簡単のために、方策パラメータ θ のうち、 i 番目の要素である θ_i のみに注意を向けることにする。そして、 $f(x; \omega_i) \equiv \omega_i^T \phi(x)$ とし、 $\nabla_{\theta_i} \ln \pi_{\theta}(x, u) \equiv \partial \ln \pi_{\theta}(x, u) / \partial \theta_i$ とし、
 $\hat{\delta}_{\lambda, K}(x, \omega_i)$ とは、 $\hat{\delta}_{\lambda, K}(x)$ の i 番目の要素であるものとする。

10

20

【 0 1 0 2 】

したがって、最小化すべき対象となる関数は、以下の式 (13) となる。

【 0 1 0 3 】

【 数 3 5 】

$$\varepsilon(\omega_i) = \frac{1}{2} \mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda, K}(x; \omega_i)^2 \mid d^{\pi}(x_{-K}) \right\} + \frac{1}{2} \mathbf{E}_{M(\theta)} \left\{ f(x; \omega_i) \right\}^2$$

30

… (13)

【 0 1 0 4 】

ここで、右辺の第 2 項は、式 (12) の制約条件のためのものである。したがって、式 (13) の偏微分は、以下ようになる。

【 0 1 0 5 】

【数 3 6】

$$\begin{aligned} \nabla_{\omega_i} \varepsilon(\omega_i) = & \mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda,K}(x; \omega_i) \nabla_{\omega_i} \hat{\delta}_{\lambda,K}(x; \omega_i) | d^n(x_{-K}) \right\} + \frac{1}{2} \nabla_{\omega_i} \mathbf{E}_{M(\theta)} \{ f(x; \omega_i) \}^2 \\ & \dots (14) \end{aligned}$$

ここで、以下の式が成り立つ。

$$\begin{aligned} \hat{\delta}_{\lambda,K}(x; \omega_i) &= \sum_{k=1}^K \lambda^{k-1} \nabla_{\theta_i} \log \pi_{\theta}(x_{-k}, u_{-k}) + \omega_i^T \nabla_{\omega_i} \hat{\delta}_{\lambda,K}(x; \omega_i) \\ \nabla_{\omega_i} \hat{\delta}_{\lambda,K}(x; \omega_i) &= (1 - \lambda) \sum_{k=1}^K \lambda^{k-1} \phi(x_{-k}) + \lambda^K \phi(x_{-K}) - \phi(x) \end{aligned}$$

従来の最小二乗法は、 $\nabla_{\omega_i} \varepsilon(\omega_i) = \mathbf{0}$ を満たすパラメータを真のパラメータ ω_i^* として見つけることを目的とするものであるが、式(13)の右辺第1項に関連する誤差 $\hat{\delta}_{\lambda,K}(x; \omega_i^*)$ とその微分 $\nabla_{\omega_i} \hat{\delta}_{\lambda,K}(x; \omega_i^*)$ との間に相関が存在すると、推定の偏りを生んでしまう。

すなわち、もしも、以下の式が成り立つとする。

$$\mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda,K}(x; \omega_i^*) \nabla_{\omega_i} \hat{\delta}_{\lambda,K}(x; \omega_i^*) | d^n(x_{-K}) \right\} \neq \mathbf{0}$$

すると、 $\nabla_{\omega_i} \varepsilon(\omega_i^*) \neq \mathbf{0}$ となる。

【0106】

一般的なRL問題では、このような相関が存在するので、このような偏りを除くために、操作変数法(instrumental variable method)を適用する[非特許文献17, 非特許文献13]。

【0107】

10

20

30

【数 3 7】

そのためには、 $\nabla_{\omega_i} \hat{\delta}_{\lambda,K}(x, \omega_i)$ は、 $\nabla_{\omega_i} \hat{\delta}_{\lambda,K}(x, \omega_i^*)$ とは相関があるものの $\hat{\delta}_{\lambda,K}(x, \omega_i^*)$ とは相関のない操作変数 $l(x)$ により置き換えられる。このような条件は、明らかに、LSTD (λ) [13, 14] と同様に、 $l(x) = \phi(x)$ のときに満たされる。式 (14) の代わりに、真のパラメータ ω_i^* 、すなわち、 $\tilde{\nabla}_{\omega_i} \varepsilon(\omega_i^*) = \mathbf{0}$ を計算するために、以下の式を 0 とするパラメータを見いだすことを試みる。

10

$$\tilde{\nabla}_{\omega_i} \varepsilon(\omega_i) =$$

$$\mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda,K}(x; \omega_i) \phi(x) \mid d^n(x_{-K}) \right\} + \mathbf{E}_{M(\theta)} \{ \phi(x) \} \mathbf{E}_{M(\theta)} \{ \phi(x) \}^T \omega_i$$

… (15)

20

【0108】

以下では、現実のマルコフ連鎖 $M(\quad)$ における時間ステップ t の状態を示すために、ノーテーションを x_t に変更する。提案する LSDG 推定アルゴリズム、LSDG (\quad) は、適格性減衰率 $\lambda \in [0, 1)$ の下で、逆方向にさかのぼる時間ステップ K を現在の状態 x_t のタイムステップ t と同じにする。すなわち、以下が成り立つ。

【0109】

【数 3 8】

$$\hat{\delta}_{\lambda,K}(x_i; \omega_i) = g_{\lambda,i}(x_i) - (z_{\lambda}(x_{i-1}) - \phi(x_i))^T \omega_i$$

ここで、

$$g_{\lambda,i}(x_i) = \sum_{s=0}^i \lambda^{i-s} \nabla_{\theta_i} \ln \pi_{\theta}(x, u)$$

であり、

$$z_{\lambda}(x_i) = (1 - \lambda) \sum_{s=1}^i \lambda^{i-s} \phi(x_s) + \lambda^i \phi(x_0) \text{ である。}$$

10

式(15)における期待値は、以下のように推定される：

$$\begin{aligned} & \lim_{K \rightarrow \infty} \mathbf{E}_{M(\theta)} \left\{ \hat{\delta}_{\lambda,K}(x; \omega_i) \phi(x) \mid d^{\pi}(x_{-K}) \right\} \\ & \cong \frac{1}{T} \sum_{t=1}^T \phi(x_t) \{ g_{\lambda,i}(x_{t-1}) - (\phi(x_t) - z_{\lambda}(x_{t-1}))^T \omega_i \} \\ & = \mathbf{b}_T - \mathbf{A}_T \omega_i \end{aligned}$$

20

ここで、以下の式が成り立つ。

$$\mathbf{b}_T \equiv \frac{1}{T} \sum_{t=1}^T \phi(x_t) g_{\lambda,i}(x_{t-1})$$

$$\mathbf{A}_T \equiv \frac{1}{T} \sum_{t=1}^T \phi(x_t) (\phi(x_t) - z_{\lambda}(x_{t-1}))^T$$

30

$$\mathbf{E}_{M(\theta)} \{ \phi(x) \} \cong \frac{1}{T+1} \sum_{t=0}^T \phi(x_t) \equiv \mathbf{c}_T$$

したがって、これらの予測子を式(15)に代入すると、時間ステップTにおける推定値 $\hat{\omega}_i^*$ は、以下のとおり計算される。

$$\mathbf{b}_T - \mathbf{A}_T \hat{\omega}_i^* + \mathbf{c}_T \mathbf{c}_T^T \hat{\omega}_i^* = 0$$

$$\Leftrightarrow \hat{\omega}_i^* = (\mathbf{A}_T - \mathbf{c}_T \mathbf{c}_T^T)^{-1} \mathbf{b}_T$$

40

【0 1 1 0】

図3は、LSDG() を求める手順をアルゴリズム1として示す図である。

【0 1 1 1】

【数 3 9】

図3では、 $\hat{\omega}_i^*$ よりも行列パラメータ Ω^* の場合についてのLSDG(λ)をアルゴリズム1として示す。

50

【 0 1 1 2 】

また、図 4 は、アルゴリズム 1 を示すフローチャートである。

図 4 を参照して、まず、ステップ 1 0 0 において、処理の前提として、以下の設定がなされる。

【 0 1 1 3 】

【 数 4 0 】

方策を $\pi_{\theta}(x, u)$ とする (パラメータ θ は固定)

状態の特徴ベクトル関数を $\phi(x)$ とする

10

【 0 1 1 4 】

続いて、初期化処理として、以下の処理が行われる (ステップ S 1 0 2)。

【 0 1 1 5 】

【 数 4 1 】

初期化 : $\lambda \in [0, 1)$, $\mathbf{c} := \mathbf{0}$; $\mathbf{z} = \mathbf{0}$; $\mathbf{g} = \mathbf{0}$; $\mathbf{A} := \mathbf{0}$; $\mathbf{B} := \mathbf{0}$

【 0 1 1 6 】

時間ステップ t が $t = 0$ とされ (ステップ S 1 0 4)、以下の処理が、 $t = 0$ から $t = T - 1$ まで繰り返される (ステップ S 1 0 6 ~ S 1 1 6)。

【 0 1 1 7 】

20

まず、ステップ S 1 0 6 において $t = 0$ であれば、初期状態が観測され (ステップ S 1 0 8)、続いて、以下の設定が行われる (S 1 1 0)。

【 0 1 1 8 】

【 数 4 2 】

$\mathbf{z} := \phi(\mathbf{x}_0)$; $\mathbf{c} := \phi(\mathbf{x}_0)$

【 0 1 1 9 】

一方、ステップ S 1 0 6 において、 t が 0 でなければ、以下の処理が行われる (ステップ S 1 1 2)。

【 0 1 2 0 】

30

【 数 4 3 】

$\mathbf{z} := \lambda \mathbf{z} + (1 - \lambda) \phi(\mathbf{x}_t)$

【 0 1 2 1 】

ステップ S 1 1 0 または S 1 1 2 に続いて、以下の計算が行われる (ステップ S 1 1 4)。

【 0 1 2 2 】

【 数 4 4 】

$\mathbf{c} := \mathbf{c} + \phi(\mathbf{x}_{t+1})$;

40

$\mathbf{g} := \lambda \mathbf{g} + \nabla_{\theta} \ln \pi_{\theta}(x_t, u_t)$;

$\mathbf{A} := \mathbf{A} + \phi(\mathbf{x}_{t+1})(\phi(\mathbf{x}_{t+1}) - \mathbf{z})^T$;

$\mathbf{B} := \mathbf{B} + \phi(\mathbf{x}_{t+1})\mathbf{g}^T$

【 0 1 2 3 】

ステップ S 1 1 6 にて、 t が T よりも小さければ処理はステップ S 1 0 6 に復帰し、 t が T 以上であれば、処理はステップ S 1 1 8 に移行して、以下の計算を行う。

【 0 1 2 4 】

50

【数 4 5】

$$\Omega := (\mathbf{A} - \mathbf{c}\mathbf{c}^T / t)^{-1} \mathbf{B};$$

【0 1 2 5】

続いて、以下の計算により LSDG の推定値を得る。

【0 1 2 6】

【数 4 6】

$$\hat{\nabla}_{\theta} \ln d^{\pi}(x) = \Omega \phi(\mathbf{x})$$

【0 1 2 7】

(4. LSDG 推定による方策の更新)

ここでは、上述した LSDG 推定に基づく PGR L アルゴリズムを定義する。

【0 1 2 8】

【数 4 7】

LSDG (λ) により $\nabla_{\theta} \ln d^{\pi}(x)$ に対する推定を実現すると、ただちに、割引率 γ とは独立に、PG に対する以下のような推定が導かれる。

$$\hat{\nabla}_{\theta} R(\theta) = \frac{1}{T} \sum_{t=0}^T (\nabla_{\theta} \ln \pi_{\theta}(x_{x_t}, u_{x_t}) + \hat{\nabla}_{\theta} \ln d^{\pi}(x_{x_t})) r_{t+1} \quad 20$$

【0 1 2 9】

方策パラメータは、適切なステップサイズ の確率的勾配法により更新される。

【0 1 3 0】

【数 4 8】

$$\theta := \theta + \alpha (\nabla_{\theta} \ln \pi_{\theta}(x_{x_t}, u_{x_t}) + \hat{\nabla}_{\theta} \ln d^{\pi}(x_{x_t})) r_{t+1}$$

【0 1 3 1】

ここで、 $:=$ は、右辺を左辺に代入することを示す。

図 5 は、LSDG () - PG を、LSDG () を利用した、PG アルゴリズムについての最も簡単な実現法の 1 つとして、アルゴリズム 2 として示す図である。ここで、減衰率パラメータ $[0, 1)$ は、古い の値により与えられる過去の推定を捨てていくために導入されている。

【0 1 3 2】

関数近似器による LSDG 推定は他の重要な内容を構成する。すなわち、特に、連続状態問題において、基底関数 $\phi(x)$ をいかにして設定するか、ということである。

【0 1 3 3】

【数 4 9】

PG 推定に対しては、LSDG に関する目的は、正確に LSDG を推定することではな

く、単に、 $\sum_{x,u} d^{\pi}(x) \pi_{\theta}(x,u) \nabla_{\theta} \ln d^{\pi}(x) r(x,u)$ を近似することである。

【0 1 3 4】

したがって、以下のような定理が有用である。

(定理 3)

【0 1 3 5】

10

20

30

40

【数50】

LSDGの推定器における基底関数を

$$\phi(x) = \sum_u \pi_\theta(x, u) r(x, u) \text{ とするとき、}$$

関数推定器 $f(x; \omega) = \omega \sum_u \pi_\theta(x, u) r(x, u)$ は、

PG $\sum_{x,u} d^\pi(x) \pi_\theta(x, u) \nabla_\theta \ln d^\pi(x) r(x, u)$ の第2項を表すことができる。 10

ここで、調整可能なパラメータ ω は、 $R^{d \times 1}$ のベクトルである：

$$\sum_{x,u} d^\pi(x) \pi_\theta(x, u) r(x, u) \nabla_\theta \ln d^\pi(x) = \sum_{x,u} d^\pi(x) \pi_\theta(x, u) r(x, u) f(x; \omega^*)$$

ここで、 ω^* は、平均誤差 $\mathcal{E}(\omega) = \frac{1}{2} \sum_x d^\pi(x) \left\{ \nabla_\theta \ln d^\pi(x) - f(x; \omega) \right\}^2$ を最 20

小化するものである。

【0136】

(証明)

以下の式により証明される。

【0137】

【数51】

$$\nabla_\omega \mathcal{E}(\omega^*) = \sum_{x,u} d^\pi(x) \pi_\theta(x, u) r(x, u) \left\{ \nabla_\theta \ln d^\pi(x) - f(x; \omega^*) \right\} = 0 \quad 30$$

【0138】

図6は、図5に対応するフローチャートである。

図6を参照して、まず、ステップ200において、処理の前提として、以下の設定がなされる。

【0139】

【数52】

方策を $\pi_\theta(x, u)$ とする (パラメータ θ は調整可能)状態の特徴ベクトル関数を $\phi(x)$ とする 40

【0140】

続いて、初期化処理として、以下の処理が行われる (ステップS202)。

【0141】

【数53】

初期化: $\lambda \in [0, 1)$, $\beta \in [0, 1)$, α ,

$$\mathbf{c} := \mathbf{0}; \mathbf{z} = \mathbf{0}; \mathbf{g} = \mathbf{0}; \mathbf{A} := \mathbf{0}; \mathbf{B} := \mathbf{0}$$

【0142】

時間ステップ t が $t = 0$ とされ (ステップS204)、以下の処理が、 $t = 0$ から $t =$ 50

T - 1 まで繰り返される (ステップ S 2 0 6 ~ S 2 1 6)。

【 0 1 4 3 】

まず、ステップ S 2 0 6 において $t = 0$ であれば、初期状態が観測され (ステップ S 2 0 8)、続いて、以下の設定が行われる (S 2 1 0)。

【 0 1 4 4 】

【数 5 4】

$$\mathbf{z} := \phi(\mathbf{x}_0); \quad \mathbf{c} := \phi(\mathbf{x}_0)$$

【 0 1 4 5 】

一方、ステップ S 2 0 6 において、 t が 0 でなければ、以下の処理が行われる (ステップ S 2 1 2)。

【 0 1 4 6 】

【数 5 5】

$$\mathbf{z} := \lambda \mathbf{z} + (1 - \lambda) \phi(\mathbf{x}_t)$$

【 0 1 4 7 】

ステップ S 2 1 0 または S 2 1 2 に続いて、以下の計算が行われる (ステップ S 2 1 4)。

【 0 1 4 8 】

【数 5 6】

$$\mathbf{c} := \beta \mathbf{c} + \phi(\mathbf{x}_{t+1});$$

$$\mathbf{g} := \beta \lambda \mathbf{g} + \nabla_{\theta} \ln \pi_{\theta}(x_t, u_t);$$

$$\mathbf{A} := \beta \mathbf{A} + \phi(\mathbf{x}_{t+1})(\phi(\mathbf{x}_{t+1}) - \mathbf{z})^T;$$

$$\mathbf{B} := \beta \mathbf{B} + \phi(\mathbf{x}_{t+1}) \mathbf{g}^T;$$

$$\Omega := (\mathbf{A} - \mathbf{c} \mathbf{c}^T / \|\mathbf{c}\|)^{-1} \mathbf{B};$$

$$\theta := \theta + \alpha \mathbf{r}_{t+1} \left\{ \nabla_{\theta} \ln \pi_{\theta}(x_t, u_t) + \Omega^T \phi(\mathbf{x}_t) \right\};$$

【 0 1 4 9 】

ステップ S 1 1 6 にて、 t が T よりも小さければ処理はステップ S 2 0 6 に復帰し、 t が T 以上であれば、処理はステップ S 2 2 0 に移行して、以下の計算を行うことで、方策がアップデートされる。

【 0 1 5 0 】

【数 5 7】

$$p(u | x; \theta) = \pi_{\theta}(x_t, u_t)$$

【 0 1 5 1 】

図 7 は、本発明の制御器の構成の概念図である。

本発明の制御器は、行動、すなわち、制御信号を制御対象に与える処理を行って、制御対象の状態量を観測器 (たとえば、位置センサ、角度センサ、加速度センサ、角加速度センサなど) で観測し、この観測結果により「定常分布の対数の偏微分」(LSDG) を推定し、方策パラメータを更新し、これにより方策を更新する。そして、更新された方策により、さらに、制御対象が制御される。

(5 . 数値計算の結果)

有限のグリッドの組 $X = \{ 1, \dots, m \}$ と 2 つの可能な行動 $U = \{ L, R \}$ (左 (L) または右 (R) への 1 グリッド分の運動) を有している「1 次元トラス状グリッド空間」において、われわれが提案したアルゴリズムのパフォーマンスを検証した。これは、典

10

20

30

40

50

型的な m 状態 MDP タスクであり、状態の遷移確率は以下のように与えられる。

【 0 1 5 2 】

【 数 5 8 】

$$p(x|x_{-1}=i, u_{-1}=L) = \begin{cases} p_i & \text{if } x=i-1 \\ \frac{1-p_i}{2} & \text{if } x=i \text{ or } i+1 \\ 0 & \text{otherwise,} \end{cases}$$

$$p(x|x_{-1}=i, u_{-1}=R) = \begin{cases} p_i & \text{if } x=i+1 \\ \frac{1-p_i}{2} & \text{if } x=i \text{ or } i-1 \\ 0 & \text{otherwise,} \end{cases}$$

10

【 0 1 5 3 】

ここで、 $x=0$ および $x=m$ ($x=1$ および $x=m+1$) とは、同じ状態であり、 p_i [0 , 1] ($i=1, \dots, m$) は、タスクに依存する定数である。われわれの数値計算では、確率的方策は、以下のようなシグモイダル関数で表される：

【 0 1 5 4 】

【 数 5 9 】

$$\pi_\theta(x, u=L) = 1 - \pi_\theta(x, u=R) = 1 / (1 + \exp(\theta^T \phi(x)))$$

【 0 1 5 5 】

ここで、状態特徴ベクトル $(1), \dots, (m)$ \mathbb{R}^m の全ての要素は、定常正規分布 $N(0, 1^2)$ からシミュレーションごとに独立に取り出された。これは、確率的方策のパラメータ化がいかにして、われわれのアルゴリズムのパフォーマンスに影響を与えるかを検証するためであった。状態特徴ベクトル (x) は、LSDG 推定のための基底関数としても使用された。各シミュレーションは、 10^5 タイムステップ以上からなる 1 つのエピソードを実行した。

30

【 0 1 5 6 】

まず、われわれは、問題設定や方策パラメータに関わりなく、LSDG () がどれくらい正確に以下の「定常分布の対数の偏微分」(LSDG) を推定しているのかを検証した。

【 0 1 5 7 】

【 数 6 0 】

$$\nabla_\theta \ln d^\pi(x)$$

【 0 1 5 8 】

この目的を達成するために、タスク依存の定数 p_1, \dots, p_m は、区間 [0 . 7 , 1] の均一な分布から独立に選ばれ、各シミュレーションでは固定された。方策パラメータ θ は、正規分布 $N(0, 0.5^2)$ に従ってランダムに選択され、各シミュレーションでは固定された。

40

【 0 1 5 9 】

図 8 は、 $m=3$ のときに、推定された LSDG の典型的な時間経過を示す図である。ここでは、9 の異なった色が、LSDG の異なった要素を示している。実線は、LSDG (0) により推定された値を示し、点線は、LSDG の解析的な解を示している。

【 0 1 6 0 】

図 8 に示すように、LSDG (0) による推定は、 $m=3$ の場合は約 100 回のタイムステップで解析的な解に収束している。

50

【0161】

7 状態タスクを用いて、適格性減衰率の影響を調べた。異なった設定についての平均的なパフォーマンスを評価するために、以下で定義される「相対誤差」基準を採用した。

【0162】

【数61】

$$E_{M(\theta)} \left\{ (f(x; \Omega^*) - f(x; \Omega))^2 \right\} / E_{M(\theta)} \left\{ (f(x; \Omega^*))^2 \right\}$$

【0163】

ここで、 θ^* は定理3において定義された最適パラメータであり、解析的に計算された。図9および図10は、 $\gamma = 0, 0.3, 0.9$ および 1 について、相対誤差の200シミュレーションについての平均の時間推移を示している。これら2つの図の間の相違は、ただ、特徴ベクトル $\phi(x)$ の要素の数だけである。

10

【0164】

図9において使用される特徴ベクトル $\phi(x) \in R^7$ は、適切なものであり、異なった状態を区別するのに十分であった一方、図10において使用される特徴ベクトル $\phi(x) \in R^6$ は、適切でない。これらの結果は、理論的な予想と合致するものである。つまり、もしも基底関数が適切であれば(図9)、 $\gamma = 1$ 以外の任意の値に θ を設定できるのに対して、そうでないときは、 $\gamma = 1$ 以外の大きな値に θ を設定する必要がある(図10)。

【0165】

20

最後に、LSDG() PGを他の従来からのPG法とを、3状態タスクにおいて比較した。ここで、状態遷移確率は、全ての $i \in \{1, 2, 3\}$ について $p_i = 1$ に設定された。

【0166】

図11は、このタスクにおける報酬の設定を示す図である。ここでは、2種類の報酬がある。定数 " $r = (\pm) 2$ " と変数 " $r = (\pm) s$ " である。変数 s は、各シミュレーションにおいて、区間 $[0.8, 1)$ での均一分布からランダムに設定された。報酬 s は、最適方策を見いだすための θ の最小値を以下のように規定していることに注意されたい：

【0167】

【数62】

$$\gamma^2 + \gamma > 2s / (2 - s)$$

30

【0168】

したがって、 θ の設定は重要であり、このタスクにおいては困難である。従来PG法のパフォーマンスのベースラインとして、2つのアルゴリズムを採用した：GPOMDP [非特許文献1] とKondaのactor-critic法とである [非特許文献5]。

【0169】

図12は、100回のシミュレーションについての3つの方法の平均のパフォーマンスを示す図である。エラーバーは100回についての標準偏差を示す。ここでは、パフォーマンスは、 $3R(\theta) / (2 - 2s)$ により評価された。すなわち、平均の報酬がその上限値である $(2 - 2s) / 3$ により正規化されている。結果はLSDG() - PG法は他の方法よりもパフォーマンスが優れていることを示している。

40

【0170】

以上説明したように、本発明では、現実の順方向および逆方向のマルコフ連鎖は密接に関連しており、定理において共通の性質を有することを利用して、これらを用いて、定常分布の対数の偏微分(LSDG)を推定するアルゴリズムとしてLSDG() を、LSDG推定を用いたPGアルゴリズムとしてLSDG() PGを提案した。実験結果はLSDG() は、適格性減衰率 $\gamma \in [0, 1)$ で動作することができ、かつ、LSDG() - PGは、割引率 γ とは独立に学習をすることができる。

【0171】

50

今回開示された実施の形態はすべての点で例示であって制限的なものではないと考えられるべきである。本発明の範囲は上記した説明ではなくて特許請求の範囲によって示され、特許請求の範囲と均等の意味および範囲内でのすべての変更が含まれることが意図される。

【図面の簡単な説明】

【0172】

【図1】本発明の制御方法および制御プログラムが適用される制御器を用いたシステム1000の一例を示す概念図である。

【図2】コンピュータ100の構成をブロック図形式で示す図である。

【図3】LSDG()を求める手順をアルゴリズム1として示す図である。

【図4】アルゴリズム1を示すフローチャートである。

【図5】LSDG()-PGを示す図である。

【図6】図5に対応するフローチャートである。

【図7】本発明の制御器の構成の概念図である。

【図8】 $m = 3$ のときに、推定されたLSDGの典型的な時間経過を示す図である。

【図9】 $m = 7$ のときに、十分な特徴ベクトルを用いて推定されたLSDGの相対誤差の時間経過を示す図である。

【図10】 $m = 7$ のときに、不十分な特徴ベクトルを用いて推定されたLSDGの相対誤差の時間経過を示す図である。

【図11】タスクにおける報酬の設定を示す図である。

【図12】100回のシミュレーションについての3つの方法の平均のパフォーマンスを示す図である。

【符号の説明】

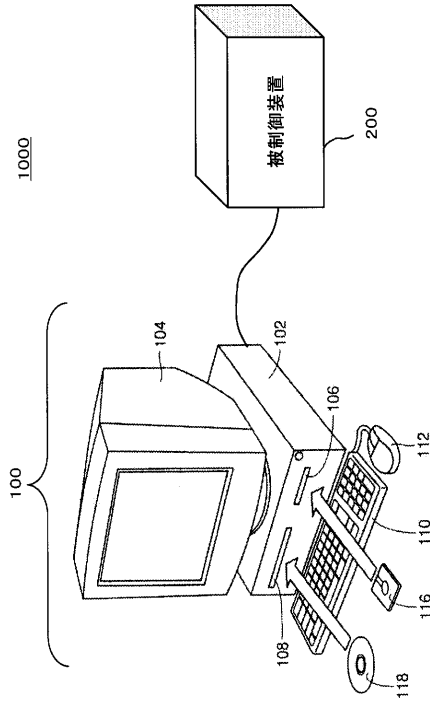
【0173】

100 コンピュータ、102 コンピュータ本体、104 ディスプレイ、106 FDドライブ、108 CD-ROMドライブ、110 キーボード、112 マウス、118 CD-ROM、120 CPU、122 メモリ、124 ハードディスク、128 通信インタフェース、200 被制御装置、1000 システム。

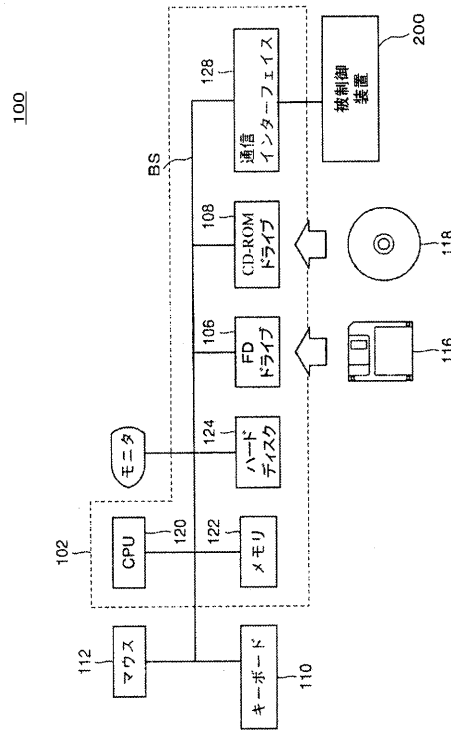
10

20

【図1】



【図2】



【図3】

アルゴリズム1
 $LSDG(\lambda): \nabla_{\theta} \ln d^{\pi}(x)$ に対する推定

所与の条件

- a policy $\pi_{\theta}(x, u)$ with a fixed θ .
- a feature vector function of state $\phi(x)$.

初期化 : $\lambda \in [0, 1)$.

Set: $c := 0; z := 0; g := 0; A := 0; B := 0$.

for $t = 0$ to $T - 1$ do

 if $t = 0$ then

$z := \phi(x_0); c := \phi(x_0);$

 else

$z := \lambda z + (1 - \lambda)\phi(x_t);$

 end if

$c := c + \phi(x_{t+1});$

$g := \lambda g + \nabla_{\theta} \ln \pi_{\theta}(x_t, u_t);$

$A := A + \phi(x_{t+1})(\phi(x_{t+1}) - z)^T;$

$B := B + \phi(x_{t+1})g^T;$

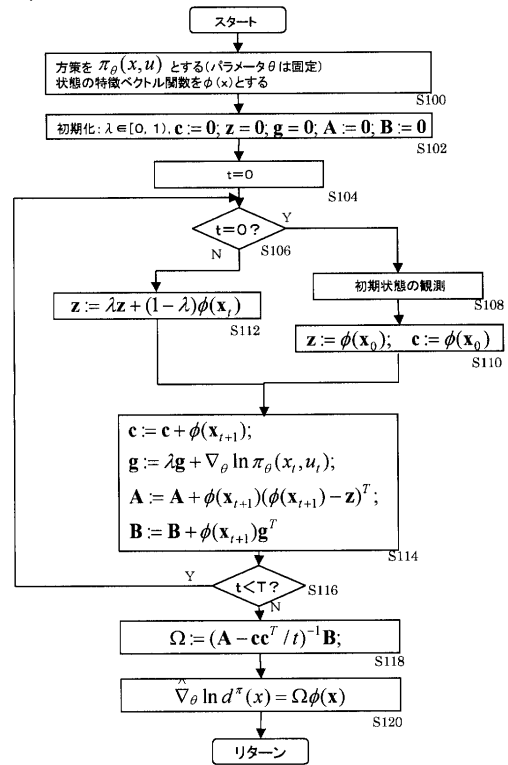
end for

$\Omega := (A - cc^T/t)^{-1}B;$

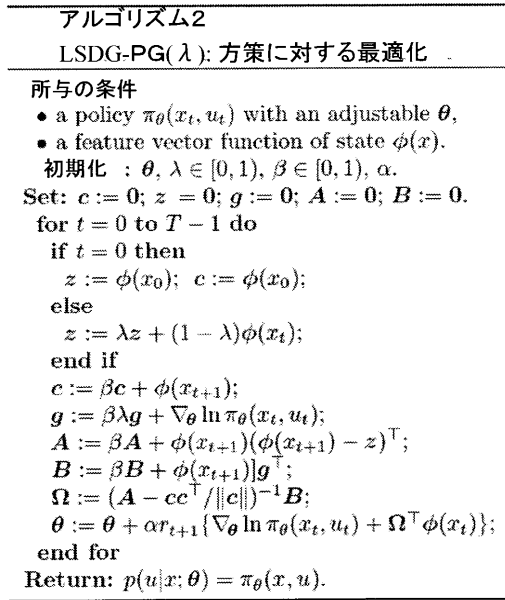
Return: $\nabla_{\theta} \ln d^{\pi}(x) = \Omega \phi(x).$

【図4】

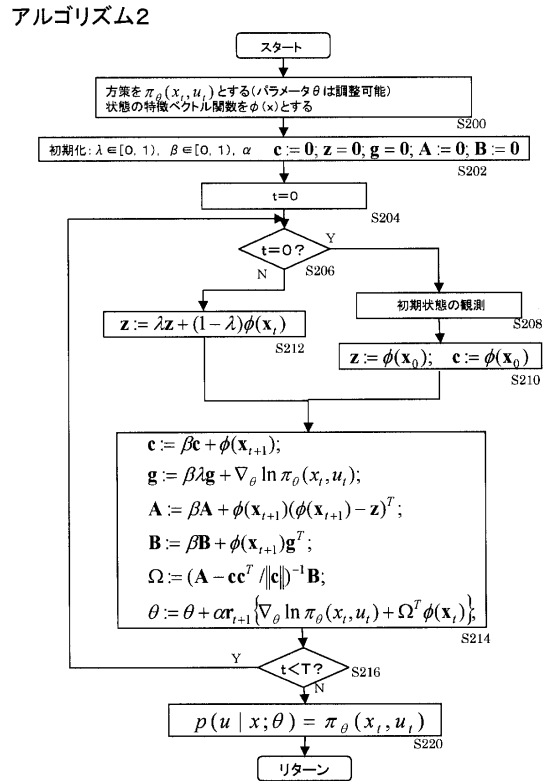
アルゴリズム1



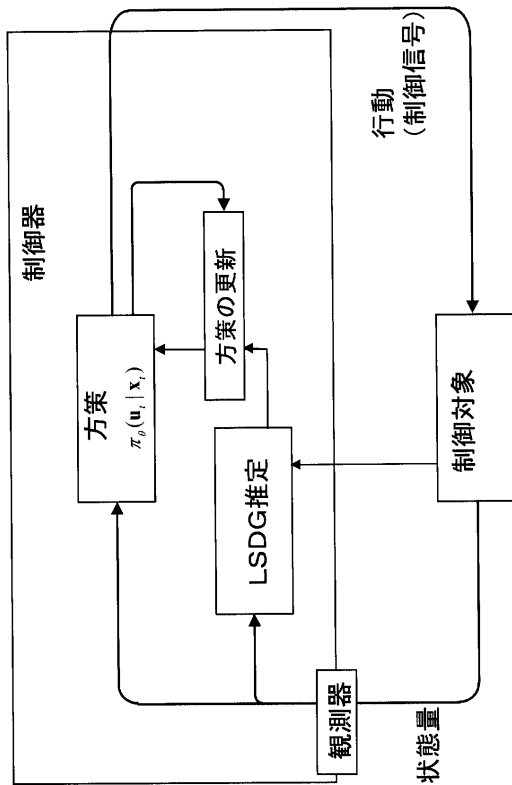
【 図 5 】



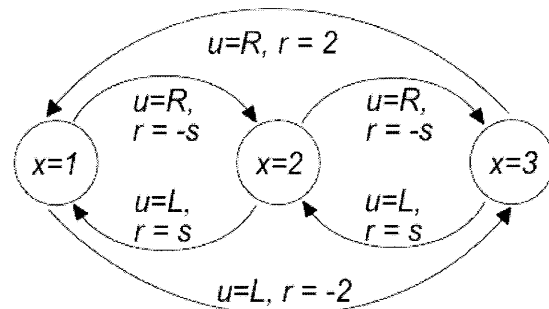
【 図 6 】



【 図 7 】

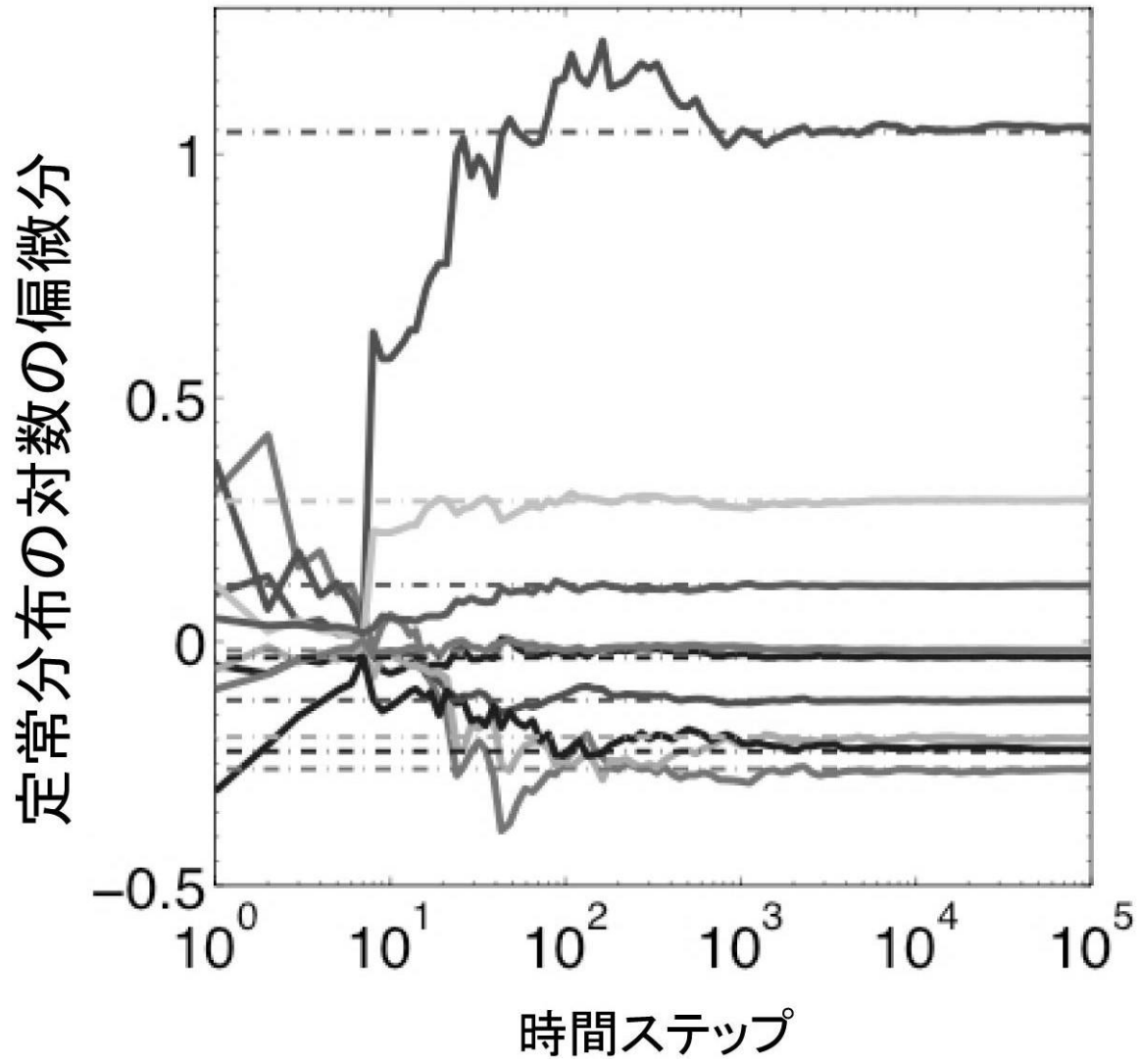


【 図 1 1 】

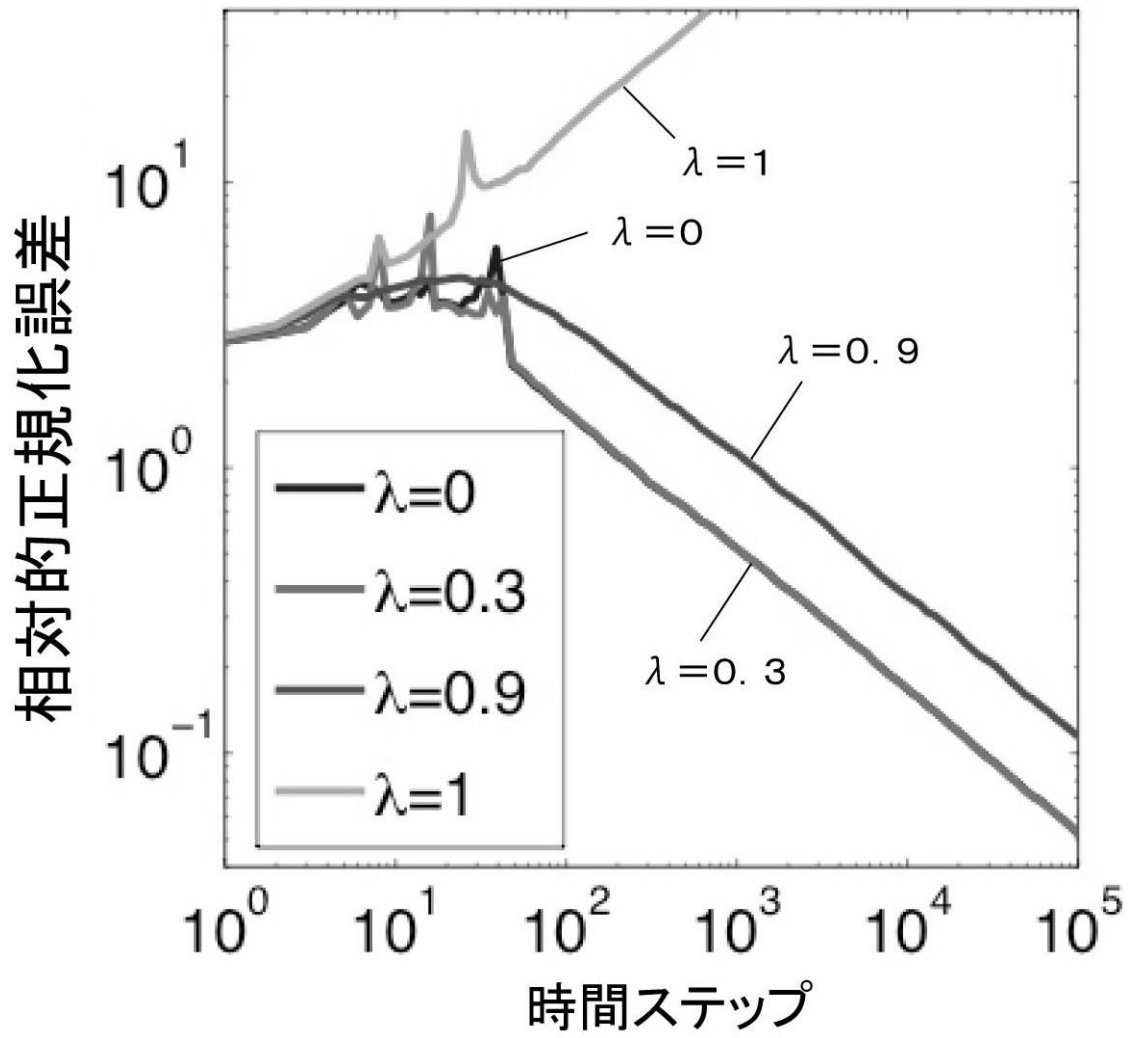


3状態マルコフ決定過程問題

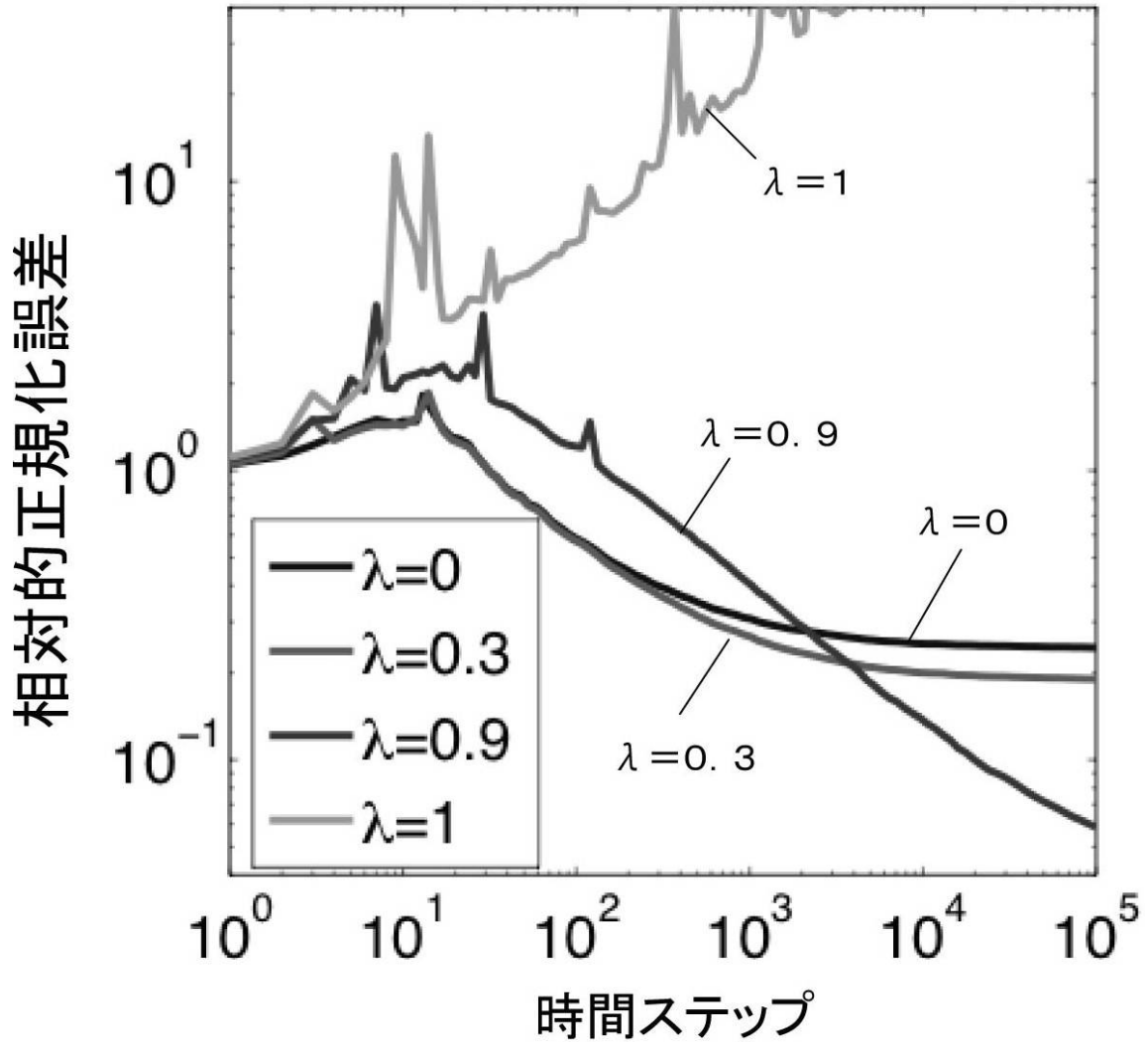
【図8】



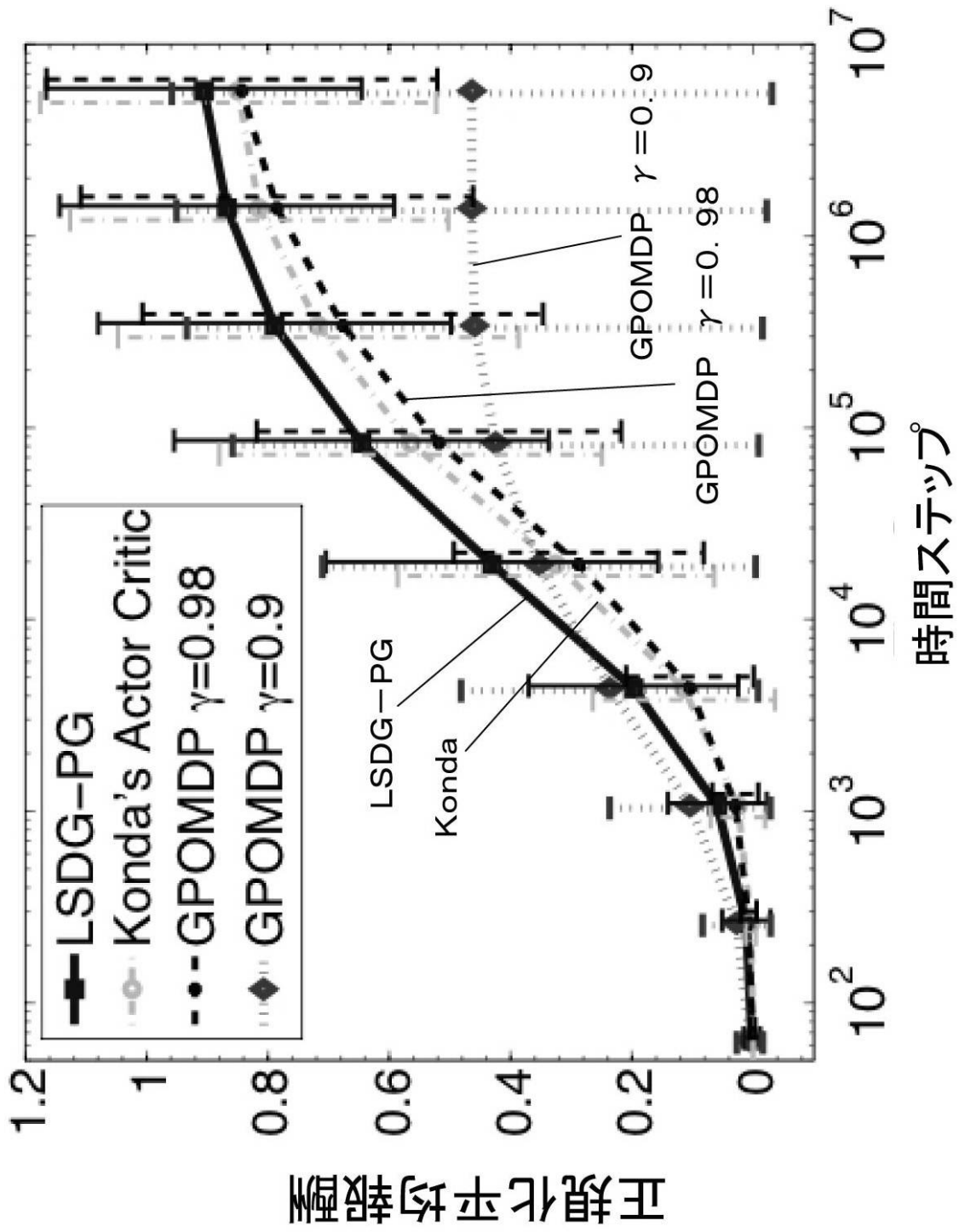
【図9】



【図10】



【 図 1 2 】



フロントページの続き

- (74)代理人 100109162
弁理士 酒井 将行
- (74)代理人 100111246
弁理士 荒川 伸夫
- (72)発明者 森村 哲郎
沖縄県国頭郡恩納村字恩納7542番地 独立行政法人沖縄科学技術研究基盤整備機構内
- (72)発明者 内部 英治
沖縄県国頭郡恩納村字恩納7542番地 独立行政法人沖縄科学技術研究基盤整備機構内
- (72)発明者 吉本 潤一郎
沖縄県国頭郡恩納村字恩納7542番地 独立行政法人沖縄科学技術研究基盤整備機構内
- (72)発明者 銅谷 賢治
沖縄県国頭郡恩納村字恩納7542番地 独立行政法人沖縄科学技術研究基盤整備機構内

審査官 柿崎 拓

- (56)参考文献 特開2007-065929(JP,A)
特開2005-084834(JP,A)

- (58)調査した分野(Int.Cl., DB名)
G05B 11/00-13/04