

(19)日本国特許庁(JP)

(12)公開特許公報(A)

(11)特許出願公開番号
特開2022-124284
(P2022-124284A)

(43)公開日 **令和4年8月25日(2022. 8. 25)**

(51)Int. Cl.	F I	テーマコード (参考)
<i>G 0 6 N 20/00 (2019. 01)</i>	G 0 6 N 20/00	5 L 0 4 9
<i>G 0 6 N 99/00 (2019. 01)</i>	G 0 6 N 99/00 1 8 0	
<i>G 0 6 Q 10/04 (2012. 01)</i>	G 0 6 Q 10/04	

審査請求 未請求 請求項の数 4 O L (全 12 頁)

(21)出願番号	特願2021-21962(P2021-21962)	(71)出願人	501241645 学校法人 工学院大学 東京都新宿区西新宿1丁目24番2号
(22)出願日	令和3年2月15日(2021. 2. 15)	(71)出願人	512155478 学校法人沖縄科学技術大学院大学学園 沖縄県国頭郡恩納村字谷茶1919番地1
		(74)代理人	110001519 特許業務法人太陽国際特許事務所
		(72)発明者	竹川 高志 東京都新宿区西新宿一丁目24番2号 学 校法人工学院大学内
		(72)発明者	高橋 春輝 東京都新宿区西新宿一丁目24番2号 学 校法人工学院大学内

最終頁に続く

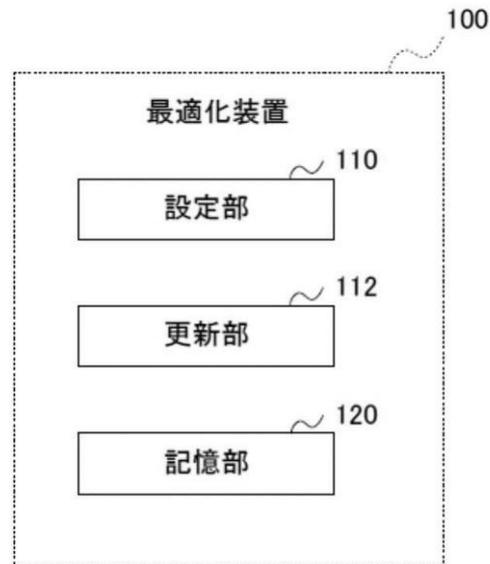
(54)【発明の名称】最適化装置、最適化方法、及び最適化プログラム

(57)【要約】

【課題】演算に係る効率を向上させることを可能とする。

【解決手段】隠れ状態を所定の態様に変更した独自隠れ状態、及び独自隠れ状態における現在の状態の推定を保持し、観測状態の法則は、条件付き確率の条件として時刻 t の観測を用い、独自隠れ状態を得るように、状態の推移法則は、条件付き確率の条件として時刻 t の独自隠れ状態及び時刻 t + 1 の独自隠れ状態を用い、エージェントの行動を得るように、状態の推移法則、及び観測状態の法則を定義する。各法則を用いて、エージェントの手順に従って分布を更新する。

【選択図】図 2



【特許請求の範囲】

【請求項 1】

状態の推移法則、観測状態の法則、及び報酬の法則による各法則が定義されている系を用い、エージェントの行動を繰り返して前記各法則を学習し報酬を獲得するモデルにおいて、

隠れ状態を所定の態様に変更した独自隠れ状態、及び前記独自隠れ状態における現在の状態の推定を保持し、

前記観測状態の法則は、条件付き確率の条件として時刻 t の観測を用い、前記独自隠れ状態を得るように、

前記状態の推移法則は、条件付き確率の条件として時刻 t の前記独自隠れ状態及び時刻 $t + 1$ の前記独自隠れ状態を用い、前記エージェントの行動を得るように、

前記状態の推移法則、及び前記観測状態の法則を定義する設定部と、

前記各法則をもとにサンプリングした確率を表す各パラメータの分布と、前記現在の状態の推定とを仮定して、ベルマン方程式に基づいて前記エージェントの最適行動を決定し、

前記各法則、所定の事前分布、及び前記最適行動を含む観測情報に対してベイズの定理を適用して得られた事後分布により、前記現在の状態の推定、及び前記各法則を用いた前記分布を更新することを繰り返す更新部と、

を含む最適化装置。

【請求項 2】

前記更新部は、エージェントの動作の手順において、各パラメータをサンプリングすることにより、ベルマン方程式に基づき、長期報酬 q が最大の行動を選択する請求項 1 に記載の最適化装置。

【請求項 3】

状態の推移法則、観測状態の法則、及び報酬の法則による各法則が定義されている系を用い、エージェントの行動を繰り返して前記各法則を学習し報酬を獲得するモデルにおいて、

隠れ状態を所定の態様に変更した独自隠れ状態、及び前記独自隠れ状態における現在の状態の推定を保持し、

前記観測状態の法則は、条件付き確率の条件として時刻 t の観測を用い、前記独自隠れ状態を得るように、

前記状態の推移法則は、条件付き確率の条件として時刻 t の前記独自隠れ状態及び時刻 $t + 1$ の前記独自隠れ状態を用い、前記エージェントの行動を得るように、

前記状態の推移法則、及び前記観測状態の法則を定義し、

前記各法則をもとにサンプリングした確率を表す各パラメータの分布と、前記現在の状態の推定とを仮定して、ベルマン方程式に基づいて前記エージェントの最適行動を決定し、

前記各法則、所定の事前分布、及び前記最適行動を含む観測情報に対してベイズの定理を適用して得られた事後分布により、前記現在の状態の推定、及び前記各法則を用いた前記分布を更新することを繰り返す、

処理をコンピュータに実行させる最適化方法。

【請求項 4】

コンピュータを、請求項 1 又は請求項 2 に記載の最適化装置の各部として機能させるためのプログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、最適化装置、最適化方法、及び最適化プログラムに関する。

【背景技術】

【0002】

10

20

30

40

50

問題設定に対する解決手法のアプローチとして、隠れマルコフモデル、ベイズ推定、及び強化学習等の手法が用いられている。

【0003】

例えば、特許文献1には、隠れ状態数および観測確率の種類と共にモデルの候補数が指数的に増加しても高速にモデル選択を実現できる隠れ変数モデル推定装置が開示されている。この隠れ変数モデル推定装置は、周辺化対数尤度関数を完全変数に対する推定量に関してラプラス近似した近似量の下界として定義される基準値を最大化することによって変分確率を計算する変分確率計算部を有する。また、隠れ変数モデル推定装置は、各隠れ状態に対して観測確率の種類とパラメータを推定することで最適な隠れ変数モデルを推定するモデル推定部と、変分確率計算部が変分確率を計算する際に用いた基準値が収束したか否かを判定する収束判定部とを有する。

10

【0004】

また、特許文献2には、環境と相互作用する強化学習エージェントが遂行する行動を選択するシステムが開示されている。このシステムは、目標回帰型ニューラルネットワーク(NN)の現在の隠れ状態に従って処理して、時間ステップについて、目標空間における初期の目標ベクトルを生成し、目標回帰型NNの内部状態を更新するように構成される、処理する工程を有している。

【0005】

また、強化学習は、状態と行動の組み合わせに対して報酬と次の状態が決定する手法である。

20

【先行技術文献】

【特許文献】

【0006】

【特許文献1】再表2013/179579号公報

【特許文献2】特表2020-508524号公報

【発明の概要】

【発明が解決しようとする課題】

【0007】

強化学習の枠組みにおいて、標準的にQ学習と呼ばれる手法が用いられている。Q学習は離散の状態に対して定義されるが、現実の課題は膨大な観測状態が存在するため、通常のQ学習では学習が難しい場合が多い。

30

【0008】

近年、発展系であるQ学習と多層のニューラルネットワークを組み合わせたDeep Q Network(DQN)がさまざまな課題において有効であることが示されている。学習済みのDQNは非常に高い性能を示すが、動作の内部状態がブラックボックスで与えられた環境をどのように解釈しているかが不明である。また、学習後の性能は高いが学習には多くの反復を必要とし、学習中に効果的に報酬を獲得することはあまり考慮されていない。

【0009】

一方、膨大な観測から重要な隠れ状態を推定しつつ状態遷移を効果的に学習するベイズ推定を用いたアルゴリズムも広く知られている。しかし、この手法では報酬の予測と状態遷移を別に扱うため、報酬と無関係な状態を詳細に分析していることとなり、問題設定によってはメモリ、及び計算量などに多大な無駄が生じる。また、多腕バンディット問題と呼ばれる枠組みにおいて、学習と報酬獲得とをバランス良く行う汎用の手法としてトンプソンサンプリングが知られているが、複雑な問題に直接適用することはできない。

40

【0010】

本発明は、上記事情を鑑みて成されたものであり、演算に係る効率を向上させることを可能とする最適化装置、最適化方法、及び最適化プログラムを提供することを目的とする。

【課題を解決するための手段】

50

【 0 0 1 1 】

上記目的を達成するために、本発明に係る最適化装置は、状態の推移法則、観測状態の法則、及び報酬の法則による各法則が定義されている系を用い、エージェントの行動を繰り返して前記各法則を学習し報酬を獲得するモデルにおいて、隠れ状態を所定の態様に変更した独自隠れ状態、及び前記独自隠れ状態における現在の状態の推定を保持し、前記観測状態の法則は、条件付き確率の条件として時刻 t の観測を用い、前記独自隠れ状態を得るように、前記状態の推移法則は、条件付き確率の条件として時刻 t の前記独自隠れ状態及び時刻 $t + 1$ の前記独自隠れ状態を用い、前記エージェントの行動を得るように、前記状態の推移法則、及び前記観測状態の法則を定義する設定部と、前記各法則をもとにサンプリングした確率を表す各パラメータの分布と、前記現在の状態の推定とを仮定して、ベルマン方程式に基づいて前記エージェントの最適行動を決定し、前記各法則、所定の事前分布、及び前記最適行動を含む観測情報に対してベイズの定理を適用して得られた事後分布により、前記現在の状態の推定、及び前記各法則を用いた前記分布を更新することを繰り返す更新部と、を含んで構成されている。

10

【 0 0 1 2 】

本発明に係る最適化方法は、状態の推移法則、観測状態の法則、及び報酬の法則による各法則が定義されている系を用い、エージェントの行動を繰り返して前記各法則を学習し報酬を獲得するモデルにおいて、隠れ状態を所定の態様に変更した独自隠れ状態、及び前記独自隠れ状態における現在の状態の推定を保持し、前記観測状態の法則は、条件付き確率の条件として時刻 t の観測を用い、前記独自隠れ状態を得るように、前記状態の推移法則は、条件付き確率の条件として時刻 t の前記独自隠れ状態及び時刻 $t + 1$ の前記独自隠れ状態を用い、前記エージェントの行動を得るように、前記状態の推移法則、及び前記観測状態の法則を定義し、前記各法則をもとにサンプリングした確率を表す各パラメータの分布と、前記現在の状態の推定とを仮定して、ベルマン方程式に基づいて前記エージェントの最適行動を決定し、前記各法則、所定の事前分布、及び前記最適行動を含む観測情報に対してベイズの定理を適用して得られた事後分布により、前記現在の状態の推定、及び前記各法則を用いた前記分布を更新することを繰り返す、処理をコンピュータに実行させる。

20

【 発明の効果 】

【 0 0 1 3 】

本発明の最適化装置、最適化方法、及び最適化プログラムによれば、演算に係る効率を向上させることを可能とする、という効果が得られる。

30

【 図面の簡単な説明 】

【 0 0 1 4 】

【 図 1 】 状態及び法則の推定に関して、従来手法の遷移図と、本実施形態の手法の遷移図との一例を示した図である。

【 図 2 】 本発明の実施形態に係る最適化装置の各機能構成を示す図である。

【 図 3 】 最適化装置のハードウェア構成を示すブロック図である。

【 図 4 】 本発明の実施形態に係る最適化装置の最適化処理ルーチンを示す図である。

【 図 5 】 本実施形態の手法と他の手法の実験結果の一例を示すグラフである。

40

【 図 6 】 実験における収束時の隠れ状態数を表にした図である。

【 発明を実施するための形態 】

【 0 0 1 5 】

以下、図面を参照して本発明の実施形態を詳細に説明する。

【 0 0 1 6 】

まず、本発明の実施形態における原理的な説明をする。

【 0 0 1 7 】

図 1 は、状態及び法則の推定に関して、従来手法の遷移図と、本実施形態の手法の遷移図との一例を示した図である。まず基本的な原理として、従来手法の状態及び法則の推定について説明する。従来手法、及び本実施形態の手法は共通して、状態の推移法則、観測

50

状態の法則、及び報酬の法則による各法則が定義されている系（遷移図）を用い、エージェントの行動を繰り返して各法則を学習し報酬を獲得する内部モデルを持つ。図1上は、従来手法の状態及び法則の推定の遷移図である。時刻 t に対して観測 o_t が得られ、行動 a_t を選択すると報酬 r_t と次の観測 o_{t+1} が得られる環境が与えられたとする。この場合に、割引率 γ に対する長期報酬 $v = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ をできるだけ大きくするような選択を行いたい。また、観測 o_t の背景には隠れ状態が存在し、行動によって隠れ状態が確率的に変化し、報酬も確率的に決定されるものとする。標準的には、隠れ状態を s_t とし、状態の推移法則 $p(s_{t+1} | s_t, a_t)$ 、観測状態の法則 $p(o_t | s_t)$ 及び報酬の法則 $p(r_t | s_t, a_t)$ が定義されている系を想定する。以下、推移法則、観測法則、及び報酬法則という。ただし、目的を達成するためのアルゴリズム（エージェント）にとってこれらの法則は未知であり、行動を繰り返して法則を学習しつつ並行して高い報酬を獲得する必要がある。

10

【0018】

一方、本発明の実施形態に係る原理において、エージェントは、環境に対する内部モデルとして、独自隠れ状態 s_t' と現在の状態の推定 $q(s_t')$ を保持している。ここで、 s_t' は s_t の報酬に関連する要素に着目して簡略化したものを想定している。ただし、 s_t は存在そのものと推移法則は仮定しているもの実際に推定するわけではない。何らかの複雑な状態と推移法則があると、それを直接考えることなく、報酬の観点から不要な状態を排除したものが s_t' である。

【0019】

20

本実施形態の手法では、推移法則、観測法則、及び報酬法則についても内部モデルを持つが、観測法則と推移法則とに関して実際の法則と異なる $p(s_t' | o_t)$ と $p(a_t | s_t', s_{t+1}')$ との形式を用いていることが特徴である。報酬法則 $p(r_t | s_t', a_t)$ に関しては実際の法則と同様である。このような形式により、本実施形態の状態及び法則の推定の遷移図は、図1下のようにできる。エージェントは各法則のパラメータを確率分布として保持し、観測結果に応じて学習する。具体的には、確率を表すパラメータである $M'(s_t', s_{t+1}', a_t) = p(a_t | s_t', s_{t+1}')$ 、 $N'(s_t', o_t) = p(o_t | s_t')$ 、 $L(s_t', a_t, r_t) = p(r_t | s_t', a_t)$ に対して、パラメータの予測である $q(M')$ 、 $q(N')$ 、 $q(L)$ が設定されている。 $q(M')$ と $q(N')$ は容易に実際の法則と対応する $q(M)$ と $q(N)$ に変換できる。

30

【0020】

ここで、パラメータについて説明する。例えば、 $p(r | s, L)$ では、 s_1 という状態で r の取り得る値が r_1, r_2, r_3 だったとすると、 $L(s_1, r_1)$ は r_1 が得られる確率を表し、 $L(s_1, r_1) + L(s_1, r_2) + L(s_1, r_3) = 1$ と、 L は行列として表現できる。 M については、 s, s', a の index をとるテンソルとなり、 a について和をとると 1 となる。

【0021】

エージェントは、次の手順で動作する。[1] 確率分布 $q(M)$ 、 $q(N)$ 、 $q(L)$ に従って M, N, L をサンプリングする。[2] サンプリングした M, N, L と状態推定 $q(s_t')$ が正しいと仮定した場合の最適行動 a_t をベルマン方程式に基づいて決定し出力する。最適行動 a_t を出力した結果、新しい情報として r_t, o_{t+1} を得る。[3] 法則 $p(a_t | s_t', s_{t+1}')$ 、 $p(o_t | s_t')$ 、 $p(r_t | s_t', a_t)$ と、事前分布 $q(s_t')$ 、 $q(M')$ 、 $q(N)$ 、 $q(L)$ と、観測された情報 o_t, a_t, r_t, o_{t+1} に対してベイズの定理を適用する。得られた事後分布を、新たな知識 $q(s_{t+1}')$ 、 $q(M')$ 、 $q(N)$ 、 $q(L)$ として更新する。[4] その後、[1] に戻り反復する。

40

【0022】

ベルマン方程式では、行動 a を行った場合の長期報酬 q の期待値が求まるので、単純に q が最大の行動を取る。ベイズの定理を適用とは、各法則 $p(a | s, s', M)$ 、 $p(o_t | s_t')$ 、 $p(r_t | s_t', a_t)$ と、事前分布 $q(s_t')$ 、 $q(M')$ 、 $q(N)$ 、 $q(L)$ と、観測された情報 o_t, a_t, r_t, o_{t+1} に対してベイズの定理を適用する。得られた事後分布を、新たな知識 $q(s_{t+1}')$ 、 $q(M')$ 、 $q(N)$ 、 $q(L)$ として更新する。

50

$s | o, L$), $p(r | s, N)$ と事前分布 $p(L, M, N)$ を用いて、事後分布 $p(s, s', L, M, N | o, o', r, r')$ を計算することを指す。以下に事後分布の計算の適用例を示す。

【0023】

$$p(s, s', L, M, N | o, o', r, r', a) \\ p(r, r', a, s, s', L, M, N | o, o') \\ = p(r | s, N) p(r' | s', N) p(a | s, s', M) p(s | o, L), p(s' | o', L) p(L, M, N)$$

【0024】

その他、計算の手法は毎回変分ベイズを用いて事後分布を収束するまで計算するが、従来手法では毎回事後分布を収束するまで計算せず変分ベイズの1stepのみ更新する実装がされている。また、本実施形態の手法の方が計算量は増えるが、オンライン性の大きな向上が見込める。

10

【0025】

上記手順の特徴について説明する。手順[1]は状態推移についてトンプソンサンプリングを適用することが従来手法では試みられていない。従来手法としては、[1]及び[2]をまとめる形でニューラルネットワークによるqの予測を行うことが主流である。また、手順[2]については、M, N, Lが既知の場合にqを求めることは標準的な手法であり、Qを元にsoftmaxで確率的に行動を決定するのが標準的な手法である。これに対して、本実施形態の手法では[1]でサンプリングしているので、softmaxは使わず単純にqが最大の行動を選択する点に特徴がある。手順[3]については、従来手法では確率モデル $p(s_{t+1} | s_t, a_t)$, $p(o_t | s_t)$ を仮定するのに対し、本実施形態の手法では、 $p(a_t | s_t', s_{t+1}')$, $p(s_t' | o_t)$ を仮定して定式化している点が大きく異なる。また、従来手法では観測 o_t を単純なカテゴリカル分布と仮定しているのに対し、本実施形態の手法ではカテゴリカル分布の直積に拡張している。

20

【0026】

なお、上記の例では、計算には変分ベイズを用いることとしているが、手順自体に変分ベイズが必須ではなく、他の計算手法を用いてもよい。

【0027】

<本発明の実施形態に係る最適化装置の構成>

次に、本発明の実施形態に係る最適化装置の構成について説明する。

30

【0028】

図2は、本発明の実施形態に係る最適化装置100の各機能構成を示す図である。図2に示すように、この最適化装置100は、機能的には、設定部110と、更新部112と、記憶部120とを備えている。

【0029】

設定部110は、状態の推移法則、及び観測状態の法則の定義を設定し、当該設定を記憶部120に保存する。以下、各法則に関して、適宜当該設定を読み出して処理を行う。

【0030】

上記原理において示したように、設定によって、隠れ状態 s_t を所定の態様に変更した独自隠れ状態 s_t' 、及び独自隠れ状態における現在の状態の推定 $q(s_t')$ を保持する。設定によって、観測状態の法則は、条件付き確率の条件として時刻 t の観測 o_t を用い、独自隠れ状態 s_t' を得るようにする ($p(s_t' | o_t)$)。設定によって、状態の推移法則は、条件付き確率の条件として時刻 t の独自隠れ状態 s_t' 及び時刻 $t+1$ の独自隠れ状態 s_{t+1}' を用い、エージェントの行動 a_t を得るようにする ($p(a_t | s_t', s_{t+1}')$)。

40

【0031】

更新部112は、エージェントの最適行動 a_t を決定し、分布を更新することを繰り返す。更新は、予め定めた条件を満たすまで繰り返せばよい。エージェントの最適行動 a_t

50

は、まず、上記エージェントの動作の手順 [1] に従って、各法則をもとに確率を表す各パラメータの M , N , L をサンプリングする。次に手順 [2] に従って、サンプリングした各パラメータ M , N , L と、現在の状態の推定 $q (s_t')$ とを仮定して、ベルマン方程式に基づいてエージェントの最適行動 a_t を決定し出力する。最適行動 a_t を出力した結果、新しい情報として r_t , o_{t+1} を得る。手順 [3] に従って、各法則 $p (a_t | s_t' , s_{t+1}')$, $p (o_t | s_t')$, $p (r_t | s_t' , a_t)$ 、所定の事前分布 $q (s_t')$, $q (M')$, $q (N)$, $q (L)$ 、及び最適行動を含む観測情報 o_t , a_t , r_t , o_{t+1} に対してベイズの定理を適用して事後分布を得る。そして、更新部 112 は、得られた事後分布により、現在の状態の推定、及び各法則を用いた分布 $q (s_{t+1}')$, $q (M')$, $q (N)$, $q (L)$ を更新することを繰り返すことにより、最終的に収束された分布を出力する。

10

【 0032 】

記憶部 120 には、設定部 110 で設定された各法則に係る設定、更新部 112 の計算過程の計算データ、及び計算結果が保存される。

【 0033 】

図 3 は、最適化装置 100 のハードウェア構成を示すブロック図である。

【 0034 】

図 3 に示すように、最適化装置 100 は、CPU (Central Processing Unit) 11、ROM (Read Only Memory) 12、RAM (Random Access Memory) 13、ストレージ 14、入力部 15、表示部 16 及び通信インタフェース (I / F) 17 を有する。各構成は、バス 19 を介して相互に通信可能に接続されている。

20

【 0035 】

CPU 11 は、中央演算処理ユニットであり、各種プログラムを実行したり、各部を制御したりする。すなわち、CPU 11 は、ROM 12 又はストレージ 14 からプログラムを読み出し、RAM 13 を作業領域としてプログラムを実行する。CPU 11 は、ROM 12 又はストレージ 14 に記憶されているプログラムに従って、上記各構成の制御及び各種の演算処理を行う。本実施形態では、ROM 12 又はストレージ 14 には、最適化プログラムが格納されている。

【 0036 】

ROM 12 は、各種プログラム及び各種データを格納する。RAM 13 は、作業領域として一時的にプログラム又はデータを記憶する。ストレージ 14 は、HDD (Hard Disk Drive) 又は SSD (Solid State Drive) 等の記憶装置により構成され、オペレーティングシステムを含む各種プログラム、及び各種データを格納する。

30

【 0037 】

入力部 15 は、マウス等のポインティングデバイス、及びキーボードを含み、各種の入力を行うために使用される。

【 0038 】

表示部 16 は、例えば、液晶ディスプレイであり、各種の情報を表示する。表示部 16 は、タッチパネル方式を採用して、入力部 15 として機能してもよい。

40

【 0039 】

通信インタフェース 17 は、端末等の他の機器と通信するためのインタフェースであり、例えば、イーサネット (登録商標)、FDDI、Wi-Fi (登録商標) 等の規格が用いられる。

【 0040 】

< 本発明の実施形態に係る最適化装置の作用 >

次に、本発明の実施形態に係る最適化装置 100 の作用について説明する。最適化装置 100 の各部として CPU 11 が、図 4 に示す最適化処理ルーチンを実行する。

【 0041 】

50

ステップS100では、CPU11が、状態の推移法則、及び観測状態の法則の定義を設定し、当該設定を記憶部120に保存する。

【0042】

ステップS102では、CPU11が、エージェントの動作の手順[1]に従って、各パラメータのM, N, Lをサンプリングする。

【0043】

ステップS104では、CPU11が、手順[2]に従って、サンプリングした各パラメータM, N, Lと、現在の状態の推定 $q(s_t')$ とを仮定して、ベルマン方程式に基づいてエージェントの最適行動 a_t を決定する。

【0044】

ステップS105では、CPU11が、最適行動 a_t を出力した結果、新しい情報として r_t, o_{t+1} を得る。

【0045】

ステップS106では、CPU11が、手順[3]に従って、各法則、所定の事前分布、及び最適行動 a_t を含む観測情報に対してベイズの定理を適用して事後分布を得る。

【0046】

ステップS108では、CPU11が、更新の条件を満たすか否かを判定する。条件を満たすと判定した場合にはステップS110へ移行し、条件を満たさないと判定した場合にはステップS102に戻って処理を繰り返す(「手順[4]」)。

【0047】

ステップS110では、最終的に得られた分布を出力し、処理を終了する。

【0048】

以上、説明した本発明の実施形態によれば、演算に係る効率を向上させることが可能である。また、技術のポイントは大きく、3つのポイントが挙げられる。1点目は確率モデルによる状態推定と意思決定問題とを統合した点、2点目は観測則と推移則とを通常の形式でなく独自の形式の法則とした点、3点目は手順におけるサンプリングの活用である。

【0049】

1点目の確率モデルによる状態推定と意思決定問題とを統合した点について説明する。これまで、カルマンフィルタなど観測からの隠れ状態の推定モデルについては様々な手法が提案されている。また、Q学習を代表として状態推移環境における意思決定問題についても多数の研究がある。しかし、現実として重要な問題であるにもかかわらず、両者を統合した問題については限定した取り組みしか行われていなかった。1点目の観点において本実施形態の技術は、状態推定と意思決定問題とを統合を手法である。

【0050】

2点目の独自の形式の法則とした点について説明する。本実施形態では、観測則と推移則とを、 $p(s_t' | o_t)$ と $p(a_t | s_t', s_{t+1}')$ の形式としたことである。この定式化により、確率モデルの上で観測 o_t が推定すべき値でなく、すでに与えられた決定事項として扱うことができる。通常の場合、観測 o_t は報酬と無関係にすべて別のものとして真の推移則全体を推定しようとする。一方、本実施形態の手法では観測 o_t ではなく、報酬が予測できれば十分となるため、真の推移則ではなく報酬の予測に必要な隠れ状態 s_t' のみにより推移則が再構成される。このことにより、隠れ状態の数が少なくなり、学習を効率良く行うことが可能である。人間の認知などにおいて、視覚及び聴覚などの膨大な観測に対して、意思決定に真に必要な状態は少数である。本法則を用いることにより、このような高度なメカニズムをシンプルなモデルで効果的に実現可能である。

【0051】

3点目の手順におけるサンプリングの活用について説明する。不確実な内部モデルに対して、現在の推定に基づいてできるだけ報酬を得ようとする活用か、将来の報酬のために情報を得る探索か、のどちらを行うかが重要な問題である。本実施形態では、原理において述べた(2)及び(3)の手順により活用を行うが、元々の推定の不確実性を(1)の

10

20

30

40

50

サンプリングで考慮しているので、最も効果的な探索を行うことができる。よって、状態推定モデルとサンプリングとによる意思決定を組み合わせたことによる効率化が図られている点が新規である。

【0052】

以上に示したように、本発明の実施形態の技術は、演算に係る効率を向上させることを可能とする。すなわち、状態と行動に依存して隠れ状態が推移し、観測状態が生成される隠れマルコフモデルに、報酬が付加されたモデルに対して、機会損失を少なくしつつ、少ない回数で状態の遷移を正しく推定し、結果の説明性が高いアルゴリズムを提供する。

【0053】

機会損失とは、エージェントの行動に関する損失である。例えば、 a_1 という行動を取れば報酬が r_1 得られるにも関わらず、 a_2 という行動で r_2 を得た場合、 $r_1 - r_2$ が機会損失となる。初期状態で情報が不完全な場合には機会損失を0にすることはできないため、不完全な情報に従って行動したり、情報が十分あるのにも関わらず探索的な行動を取ったりする場合に、機会損失が大きくなるという性質がある。アルゴリズムを長期間繰り返し実行した場合に平均機会損失が少ないことが重要であり、機会損失を少なくすること、トータルの報酬獲得を大きくすることとはほぼ同じ意味を表す。

【0054】

[実験結果]

本実施形態の手法の実験結果を説明する。図5は、本実施形態の手法と他の手法の実験結果の一例を示すグラフである。図5は、推定手法ごとの試行回数ごとの累積報酬を示している。本実施形態の提案手法は、 $O_{mul} S R$ である。本実施形態の提案手法が最も早く最適な報酬を得られている。また、図6は、実験における収束時の隠れ状態数を表にした図である。本実施形態の手法では、観測状態が32種類に対して報酬と関係のある8種類の隠れ状態を推定することができる。このことが、高速に学習を行える要因である。

【0055】

なお、本実施形態の法則に関して補足する。図1で示されている法則は、エージェントには知らされない「真の法則」であり、この法則においては隠れ状態から観測が生成される。通常、確率モデルにおいては真の法則を求めることが一般的であるが、本実施形態で想定している問題設定では必ずしも複雑な真の法則すべてを推定することが必要ではない。特に状態 s から観測 o が得られるという部分については、次回どのような観測が得られるかという法則を学習することになるが、実際には状態 o は観測できるため予測は不要であり、隠れ状態 s と報酬 r さえ予測できればよい。よって、 o をあえて予測しなくても定式化を考えた結果が本実施形態の手法である。

【0056】

例えば、真の法則において、4つの状態 $\{s_0, s_1, s_2, s_3\}$ と1対1で観測 $\{o_0, o_1, o_2, o_3\}$ が対応しているが、報酬の観点からは4つの状態が等価である場合を考えると、報酬を得るという目的のためにはどの観測であっても等価といえる。従来手法では、 s_0 という状態から o_0 が、 s_1 という状態から o_2 が生成されるということを区別して学習し、それに伴って状態数も多く必要となっていた。本実施形態の手法では、 o_0, o_1, o_2, o_3 が観測された場合は共通の状態であるということを学習している。真の法則というのも必ずしも1つの見方に固定されるものでなく、状態は同じで観測が確率的に生成されるという解釈もできる。その意味においてこの例では s と s' とは完全に等価に対応している。

【0057】

一般に、従来手法のエージェントは必ず観測 o が説明できるような状態とその推移モデルを構築している。一方、本実施形態の手法では様々な解釈があり得る中で、一見多様で複雑な観測に捉われずに、状態を報酬の観点からシンプルに再構成する。 $p(s|o)$ 、 $p(o)$ からベイズの定理により $p(o|s)$ が求まることから、 s と s' とは等価なモデルといえる。

10

20

30

40

50

【 0 0 5 8 】

本発明の実施形態の技術は、観測状態と報酬が与えられる広範囲の問題設定に対し適用可能であり、従来の手法に比べて高速に学習可能であるため、多岐にわたる応用が考えられる。特にロボットの分野、及びAIを用いたエージェントシステムの高性能化などに活用が期待される。

【 0 0 5 9 】

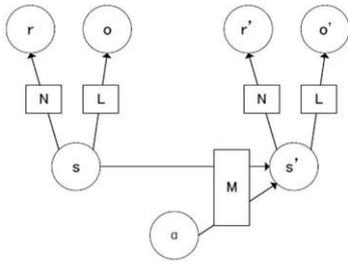
なお、本発明は、上述した実施形態に限定されるものではなく、この発明の要旨を逸脱しない範囲内で様々な変形や応用が可能である。

【符号の説明】

【 0 0 6 0 】

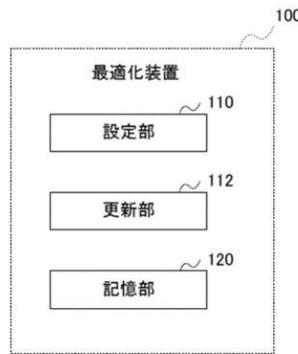
- 1 0 0 最適化装置
- 1 1 0 設定部
- 1 1 2 更新部
- 1 2 0 記憶部

【 図 1 】

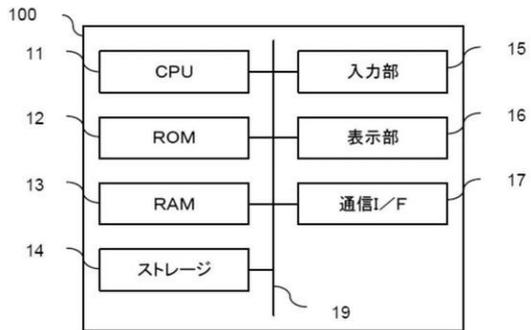


従来の状態及び法則の推定に係る遷移図

【 図 2 】

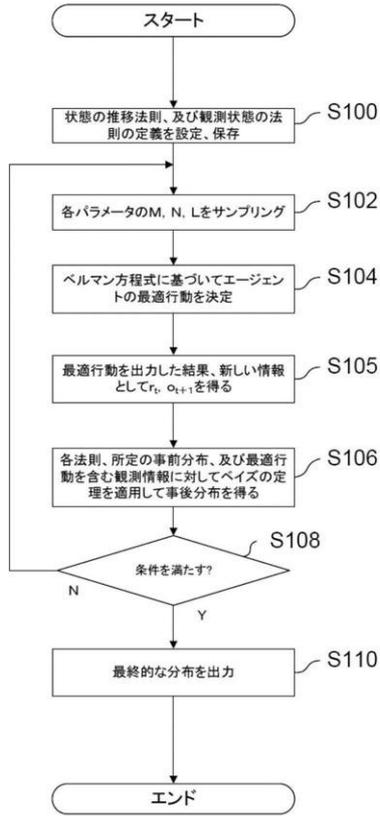


【 図 3 】

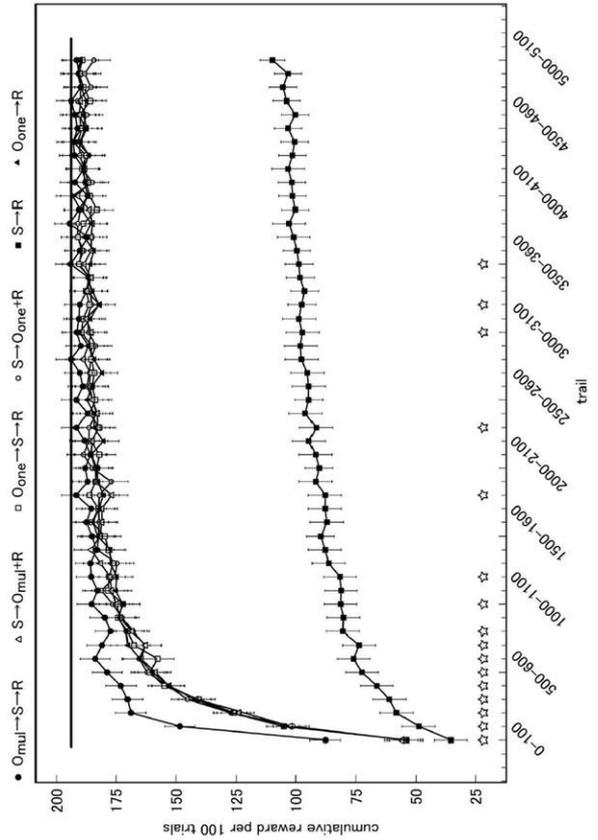


本実施形態の状態及び法則の推定に係る遷移図

【 図 4 】



【 図 5 】



【 図 6 】

モデル	O _{mul} →S→R (本提案手法)	O _{one} →S→R	S→O _{mul} +R	S→O _{one} +S
初期値 S =8	1453	1505	2543	3008
収束時の隠れ状態数	S =8	S =8	S =8	S =8
初期値 : S =32	1461	1492	3102	3078
収束時の隠れ状態数	S =8	S =32	S =32	S =32

フロントページの続き

(72)発明者 酒井 裕

東京都新宿区市谷柳町15-1-1204

(72)発明者 深井 朋樹

沖縄県国頭郡恩納村字谷茶1919番地1 学校法人沖縄科学技術大学院大学学園内

Fターム(参考) 5L049 AA04