

(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2013-205890

(P2013-205890A)

(43) 公開日 平成25年10月7日(2013.10.7)

(51) Int.Cl.	F I	テーマコード (参考)
GO6N 3/00 (2006.01)	GO6N 3/00 560A	
GO6N 5/04 (2006.01)	GO6N 5/04 550J	
GO6F 17/30 (2006.01)	GO6F 17/30 210D	

審査請求 未請求 請求項の数 6 O L (全 15 頁)

<p>(21) 出願番号 特願2012-71205 (P2012-71205)</p> <p>(22) 出願日 平成24年3月27日 (2012. 3. 27)</p> <p>特許法第30条第1項適用申請有り 平成23年11月21日に公益社団法人 計測自動制御学会 システム・情報部門の学術講演会2011 講演論文集にて発表</p>	<p>(71) 出願人 504136568 国立大学法人広島大学 広島県東広島市鏡山1丁目3番2号</p> <p>(74) 代理人 110001427 特許業務法人前田特許事務所</p> <p>(72) 発明者 保田 俊行 広島県東広島市鏡山一丁目4番1号 国立大学法人広島大学大学院工学研究院内</p> <p>(72) 発明者 大倉 和博 広島県東広島市鏡山一丁目4番1号 国立大学法人広島大学大学院工学研究院内</p>
--	---

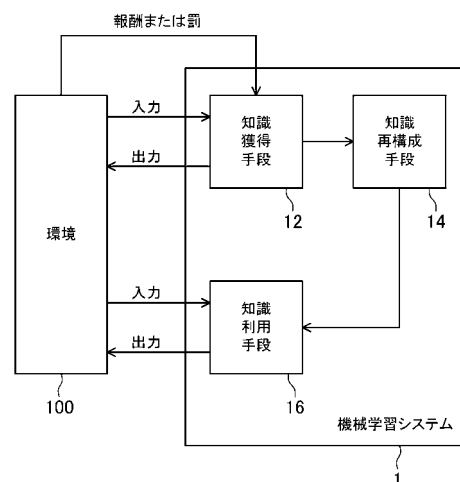
(54) 【発明の名称】 機械学習システムおよび機械学習方法

(57) 【要約】

【課題】パラメトリック表現された状態空間とノンパラメトリック表現された状態空間とを適応的に選択する。

【解決手段】機械学習システム(1)は、入力および当該入力に対する出力に対して与えられる報酬または罰に基づいて強化学習を行って、パラメトリック表現されたクラス集合を生成する知識獲得手段(12)と、パラメトリック表現されたクラス集合の生成に使用された学習済み入力に基づいて、ノンパラメトリック表現されたクラス集合を生成する知識再構成手段(14)と、未知の入力がノンパラメトリック表現されたどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をする知識利用手段(16)とを備えている。知識再構成手段(14)は、学習済み入力の個数が所定数よりも多く、かつ、パラメトリック表現された各クラスの分散が所定値よりも小さいとき、ノンパラメトリック表現されたクラス集合を生成する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

入力が状態空間におけるどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をし、入出力を繰り返すことで環境に適応した知識を獲得する機械学習システムであって、

入力および当該入力に対する出力に対して与えられる報酬または罰に基づいて強化学習を行って、パラメトリック表現されたクラス集合を生成する知識獲得手段と、

前記パラメトリック表現されたクラス集合の生成に使用された学習済み入力に基づいて、ノンパラメトリック表現されたクラス集合を生成する知識再構成手段と、

未知の入力が前記ノンパラメトリック表現されたどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をする知識利用手段とを備え、

前記知識再構成手段は、前記学習済み入力の個数が所定数よりも多く、かつ、前記パラメトリック表現された各クラスの分散が所定値よりも小さいとき、前記ノンパラメトリック表現されたクラス集合を生成する

ことを特徴とする機械学習システム。

【請求項 2】

請求項 1 に記載の機械学習システムにおいて、

前記パラメトリック表現された各クラスが多変量の正規確率分布であり、

前記知識獲得手段は、ベイズ判別法に従って、入力が前記パラメトリック表現されたどのクラスに属するかクラス判別を行う

ことを特徴とする機械学習システム。

【請求項 3】

請求項 1 および 2 のいずれか一つに記載の機械学習システムにおいて、

前記知識再構成手段は、サポートベクターマシンを用いて前記学習済み入力を線形分離して、前記ノンパラメトリック表現されたクラス集合を生成する

ことを特徴とする機械学習システム。

【請求項 4】

入力が状態空間におけるどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をし、入出力を繰り返すことで環境に適応した知識を獲得する機械学習方法であって、

入力および当該入力に対する出力に対して与えられる報酬または罰に基づいて強化学習を行って、パラメトリック表現されたクラス集合を生成する第 1 のステップと、

前記パラメトリック表現されたクラス集合の生成に使用された学習済み入力に基づいて、ノンパラメトリック表現されたクラス集合を生成する第 2 のステップと、

未知の入力が前記ノンパラメトリック表現されたどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をする第 3 のステップとを備え、

前記第 2 のステップでは、前記学習済み入力の個数が所定数よりも多く、かつ、前記パラメトリック表現された各クラスの分散が所定値よりも小さいとき、前記ノンパラメトリック表現されたクラス集合が生成される

ことを特徴とする機械学習方法。

【請求項 5】

請求項 4 に記載の機械学習方法において、

前記パラメトリック表現されたクラスが多変量の正規確率分布であり、

前記第 1 のステップでは、ベイズ判別法に従って、入力が前記パラメトリック表現されたどのクラスに属するかクラス判別が行われる

ことを特徴とする機械学習方法。

【請求項 6】

請求項 4 および 5 のいずれか一つに記載の機械学習方法において、

前記第 2 のステップでは、サポートベクターマシンを用いて前記学習済み入力が線形分離され、前記ノンパラメトリック表現されたクラス集合が生成される

10

20

30

40

50

ことを特徴とする機械学習方法。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、機械学習に関し、特に、強化学習による知識獲得の頑健性向上に関する。

【背景技術】

【0002】

システムを制御する場合、一般的にはモデル化に基づくトップダウン的アプローチがとられる。しかし、システムの大規模化などの要因により制御が困難になるということも考えられる。一方、ボトムアップアプローチではシステムの構成要素を智能化することで系全体としての合目的な入出力関係の獲得が可能である。その中の一つに強化学習法がある。強化学習法は、目標状態を与えるのみでそこに至る入出力の系列を自律的に構築できるという実装の容易さから、さまざまなシステムへの応用が期待される。

【0003】

強化学習法の従来のはず組みでは、離散的な状態・行動空間における写像関係の構築を対象としている。ここで、学習性能はこの状態・行動空間の離散化具合に大きく影響されるが、現在のところそのための設計指針は存在していない。この問題は、連続空間において動作する多くの実システムでは重大な課題である。本願発明者らはこの状態・行動空間の設計問題に対する手法として、強化学習を機能拡張したBayesian-discrimination-function-based Reinforcement Learning (BRL) を研究・開発してきた。BRLは、連続な状態・行動空間を自律的に分割する機能を持つ。さらには、従来型強化学習はマルコフ環境において学習収束が保証されているのみであるが、BRLは学習過程で分割具合を適応的に更新可能であるために動的環境でも学習可能であるという特徴を持つ。これまで、本願発明者らは、実システムとしてロボット、特に複数のロボットで構成されるマルチロボットシステム (Multi-Robot Systems: MRS) を取り上げ、自律移動ロボット群やアーム型ロボット群による協調問題においてBRLの有効性を示してきた。

【0004】

ところが、その後の追加実験において、BRLでは行動獲得後にさらに学習を続けると、徐々にその頑健性が損なわれる場合があることが観察された。これは、タスク達成に寄与しないルールは削除され、寄与するルールのみが強化されてルール集合に残ることが原因である。すなわち、BRLでは環境に特化したルール集合となる結果、過学習状態となるためにシステムが不安定になる。そこで、近年、本願発明者らはパターン認識手法の一つであるSupport Vector Machine (SVM) の高い識別性能に着目し、SVMによるルール判別がBRLの過学習抑制に有効であることを明らかにした (例えば、非特許文献1参照)。

【先行技術文献】

【非特許文献】

【0005】

【非特許文献1】J. Sakanoue, T. Yasuda, and K. Ohkura, "Preservation and Application of Acquired Knowledge Using Instance-Based Reinforcement Learning," Joint 5th International Conference on Soft Computing and Intelligent Systems and 10th International Symposium on advanced Intelligent Systems, 2010, pp.576-581

【発明の概要】

【発明が解決しようとする課題】

【0006】

BRLの過学習抑制にSVMが有効であることは実証できたものの、具体的にBRLのルール判別にどのようにSVMを用いるかについてはまだ提案できていない。かかる問題に鑑み、本発明は、機械学習システムにおいて、パラメトリック表現された状態空間とノンパラメトリック表現された状態空間とを適応的に選択する手法を提供することを目的とする。

10

20

30

40

50

【課題を解決するための手段】

【0007】

本発明の一局面に従った機械学習システムは、入力の状態空間におけるどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をし、入出力を繰り返すことで環境に適応した知識を獲得する機械学習システムであって、入力および当該入力に対する出力に対して与えられる報酬または罰に基づいて強化学習を行って、パラメトリック表現されたクラス集合を生成する知識獲得手段と、前記パラメトリック表現されたクラス集合の生成に使用された学習済み入力に基づいて、ノンパラメトリック表現されたクラス集合を生成する知識再構成手段と、未知の入力が前記ノンパラメトリック表現されたどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をする知識利用手段とを備えている。前記知識再構成手段は、前記学習済み入力の個数が所定数よりも多く、かつ、前記パラメトリック表現された各クラスの分散が所定値よりも小さいとき、前記ノンパラメトリック表現されたクラス集合を生成する。

10

【0008】

これによると、知識獲得手段による強化学習が十分に進んだところで、知識獲得手段において生成されたパラメトリック表現されたクラス集合が知識再構成手段によってノンパラメトリック表現されたクラス集合に再構成され、当該ノンパラメトリック表現されたクラス集合を用いて知識利用手段によって未知の入力のクラス判別が行われる。

【0009】

具体的には、前記パラメトリック表現された各クラスが多変量の正規確率分布であり、前記知識獲得手段は、ベイズ判別法に従って、入力が前記パラメトリック表現されたどのクラスに属するかクラス判別を行う。

20

【0010】

また、具体的には、前記知識再構成手段は、SVMを用いて前記学習済み入力を線形分離して、前記ノンパラメトリック表現されたクラス集合を生成する。

【発明の効果】

【0011】

本発明によると、機械学習システムにおいて、パラメトリック表現された状態空間とノンパラメトリック表現された状態空間とが適応的に選択され、機械学習システムの頑健性が向上する。これにより、機械学習システムが環境変化にも柔軟に対応することができるようになる。

30

【図面の簡単な説明】

【0012】

【図1】本発明の一実施形態に係る機械学習システムの機能ブロック図

【図2】図1の機械学習システムによる知識獲得および利用のフローチャート

【図3】計算機実験の実験環境を示す模式図

【図4】各SVMによるタスク達成率を示すグラフ

【図5】知識獲得までに要したエピソード数を示すグラフ

【発明を実施するための形態】

【0013】

以下、図面を参照しながら本発明を実施するための形態について説明する。なお、本発明は、以下の実施形態に限定されるものではない。

40

【0014】

(機械学習システムの実施形態)

図1は、本発明の一実施形態に係る機械学習システムの機能ブロック図である。本実施形態に係る機械学習システム1は、入力の状態空間におけるどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をし、入出力を繰り返すことで環境100に適応した知識を獲得するものである。機械学習システム1は、例えばMRSを構成する各ロボットなどに適用可能である。

【0015】

50

機械学習システム 1 は、知識獲得手段 1 2、知識再構成手段 1 4、および知識利用手段 1 6 を備えている。これら各手段は電子デバイスなどのハードウェアとして実現してもよいし、コンピュータで実行されるソフトウェアモジュールとして実現することもできる。以下、各手段について詳細に説明する。

【 0 0 1 6 】

(知識獲得手段 1 2 の詳細説明)

知識獲得手段 1 2 は、環境 1 0 0 からの入力および当該入力に対する出力に対して与えられる報酬または罰に基づいて強化学習を行って、パラメトリック表現されたクラス集合を生成する。例えば、知識獲得手段 1 2 は B R L によって強化学習を行う。すなわち、パラメトリック表現された各クラスは多変量の正規確率分布であり、知識獲得手段 1 2 は、

10

【 0 0 1 7 】

B R L

B R L では統計的にパターン分類を行うベイズ判別法を用いて入力 x が k 番目 (ただし、 k は 1 から N までの整数である。) のクラス C_k に分類されるかを識別する。ベイズ判別法は、識別対象のクラス $C = \{ C_k \}_{k=1}^K$ および各クラスの事前確率 $P(C_k)$ と確率分布 $p(x | C_k)$ が既知の場合、入力 x が観測されたときの各クラス C_k の事後確率 $P(C_k | x)$ をベイズの公式から求め、事後確率最大となるクラスに入力を識別する方法である (数式 (1) 参照)。B R L では、(1) クラスの追加と削除、(2) 確率分布モデルのパラメータ更新によって観測データから環境の確率モデルをリアルタイムに更新し、状態空間の分割を行う。

20

【 0 0 1 8 】

【 数 1 】

$$P(C_k|x) = \frac{P(C_k)p(x|C_k)}{\sum_{k=1}^K P(C_k)p(x|C_k)} \quad (1)$$

【 0 0 1 9 】

ルール構成

30

各クラスをガウス分布によって表現し、各クラスの確率分布を表すパラメータとそのときの出力を if-then 形式で記述したルールとして知識獲得手段 1 2 に記憶する。これ以降、クラスとルールを同義として扱う。ルール集合 R はルール $r_1 \dots R$ により構成され、各ルールは次式で記述される。

【 0 0 2 0 】

【 数 2 】

$$rl := \langle v, \sum, f, u, \Phi, a \rangle \quad (2)$$

40

【 0 0 2 1 】

各ルール r_1 は特徴ベクトル $v = \{ v_1, \dots, v_{n_d} \}^T$ 、共分散行列、クラスの事前確率 f 、クラスの信頼性を表す有効度 u 、各クラスで観測されたセンサ入力の集合 $\Phi = \{ \phi_1, \dots, \phi_{n_s} \}^T$ 、そして、動作 $a = \{ a_1, \dots, a_{n_a} \}^T$ より構成されている。ただし、 n_d は入力空間の次元数、 n_a は出力空間の次元数、 n_s は各クラスが記憶しているサンプルデータを表す。学習初期、状態空間にはクラスは存在せず、機械学習システム 1 が実際に観測した入出力をもとに状態空間にクラスを追加し、状態空間をガウス分布で覆っていく。

【 0 0 2 2 】

動作選択

50

以下に B R L の行動選択の概要を示す。

【 0 0 2 3 】

・ 知覚した入力を、ベイズ判別法によりどのクラスに属するかの判別を行う。

【 0 0 2 4 】

・ 既存のルールに属さない場合、ランダム行動を出力し、罰を受けなければ新たなルールを作成する。

【 0 0 2 5 】

・ 既存のルールに属す場合、そのルール行動を出力する。

【 0 0 2 6 】

入力に対する各ルールの事後確率をベイズの公式から求め、事後確率最大のルールに記述されている出力を実行する。ここでは、まず事後確率の負の対数を取り、誤って識別する確率 g_i が最小となるルールを勝者ルール r_{l_w} とする。

【 0 0 2 7 】

【 数 3 】

$$g_w = \min_i \{g_i\} \quad i \in [0, n_{r_l}] \quad (3)$$

【 0 0 2 8 】

【 数 4 】

$$g_i = -\log\{f_i \cdot p(X|C_i)\} \quad (4)$$

【 0 0 2 9 】

このとき、事後確率が非常に小さいルールが選択されるのは適切でないと考え、事後確率に閾値 P_{t_h} を設ける。そして、それをもとに計算される閾値 $g_{t_h} = -\log\{f_0 \cdot P_{t_h}\}$ によって r_{l_w} の動作を実行するかどうか判断する。具体的には、 $g_w < g_{t_h}$ の場合、 r_{l_w} の動作 A_w を実行する。 $g_w \geq g_{t_h}$ の場合、ランダムに動作を実行する。なお、 f_0 および P_{t_h} は定数である。

【 0 0 3 0 】

有効度の更新

Profit Sharing と Bucket Brigade 的戦略により報酬を過去に遡って伝播させる。その他、ループ行動を防ぐために選択されたルールに課すコスト、報酬獲得に寄与しないルールを削除してメモリ消費量を抑えるためにタスク達成時に全ルールに作用させる消散がある。

【 0 0 3 1 】

パラメータの更新

各ルールは入力をもとに確率分布のパラメータをオンラインで更新していく。リアルタイムに更新することで環境やシステム変動に対する迅速な対応が期待できる。その反面、ノイズや一時的な入力の偏りに影響を受けやすいため何らかの対処が必要となる。B R L では、区間推定法を用いたパラメータ更新によりこの問題を解決する。区間推定法は、確率分布のパラメータがある区間に入る確率を設定した確率以上になるように保証する手法であり、サンプルデータが増大するにつれて推定精度が上がる。そのため、観測データの増加に伴ってより信頼性の高いパラメータ推定が期待できる。

【 0 0 3 2 】

【 数 5 】

$$v_j \leftarrow v_j + \alpha(\bar{x}_j - v_j) \quad (5)$$

10

20

30

40

50

【 0 0 3 3 】

【 数 6 】

$$\sigma_j^2 \leftarrow \sigma_j^2 + \alpha^2 (s_j^2 - \sigma_j^2) \quad (6)$$

【 0 0 3 4 】

【 数 7 】

$$f \leftarrow \begin{cases} f + \beta(1-f) & (\text{if } P \geq 0) \\ (1-\beta)f & (\text{otherwise}) \end{cases} \quad (7)$$

10

【 0 0 3 5 】

ここで、 v は $r l_w$ の平均、 s^2 は $r l_w$ の分散、 β は区間 $[0, 1]$ の定数、 j は入力ベクトルにおける各次元、 x パーはサンプル入力の平均、 s^2 はサンプル入力の分散、 P は報酬である。

【 0 0 3 6 】

既存ルールのパラメータに基づく行動空間の適応的探索

常に行動空間をランダムに探索するのは非効率であるという観点から、知識獲得がある程度行われた状況では幅広く行動空間を探索するよりも既存のルールの近傍を探索して行動の調整を行うことが有効であると考えられる。そこで、ランダム探索をするための閾値 P_{th} の他に新たに閾値 P'_{th} を設定 ($P'_{th} < P_{th}$) し、 $g_{th} < g_w < g'_{th}$ の場合はその間にあるルールのパラメータを参照して新しいルールパラメータを決定する。つまり、行動選択を以下のように変更する。

20

【 0 0 3 7 】

・ $g_w < g_{th}$ の場合、 $r l_w$ の動作 A_w を実行する。

【 0 0 3 8 】

・ $g_{th} < g_w < g'_{th}$ の場合、この間にあるルールをもとに動作を生成する。

【 0 0 3 9 】

・ $g_w < g'_{th}$ の場合、ランダムに動作を実行する。

【 0 0 4 0 】

$g_{th} < g_w < g'_{th}$ の行動

この範囲には、 $r l_w$ 以外にも複数のルールが含まれる場合がある。これらのルールはその状況化での選択確率としては大きな差はないものの、それまでの学習過程におけるタスク達成への貢献度に従って有効度が異なる。そのため、新しいルールの動作 A' は、この範囲に含まれるルールの有効度に基づく加重平均により求める。

30

【 0 0 4 1 】

【 数 8 】

$$A' = \sum_{l=1}^{n_r} \left(\frac{u_l}{\sum_{k=1}^{n_r} u_k} \cdot A_l \right) + N(0, \sigma_{act}) \quad (8)$$

40

【 0 0 4 2 】

n_r はこの範囲に含まれるルール数であり、 $N(0, \sigma_{act})$ は平均 0・標準偏差 σ_{act} の正規分布を用いたノイズである。ノイズを付加することで、 $r l_w$ 以外のルールがない場合であっても、 A_w の近傍を探索することができる。

【 0 0 4 3 】

(知識再構成手段 1 4 および知識利用手段 1 6 の詳細説明)

知識再構成手段 1 4 は、知識獲得手段 1 2 においてパラメトリック表現されたクラス集合の生成に使用された学習済み入力に基づいて、ノンパラメトリック表現されたクラス集

50

合を生成する。例えば、知識再構成手段 14 は SVM を用いて、知識獲得手段 12 における学習済み入力を線形分離して、ノンパラメトリック表現されたクラス集合を生成する。知識利用手段 16 は、環境 100 からの未知の入力がノンパラメトリック表現されたどのクラスに属するかクラス判別を行って当該判別結果に応じた出力をする。

【0044】

SVM による知識再構成と利用

SVM は高い識別性能を持つとさまざまな分野で示されている。そこで、BR L の獲得知識のデータを SVM により識別することで、より正確な行動決定が可能になると期待できる。

【0045】

本願発明者らは、BR L によりタスク達成に有効なルールを獲得できずサンプルデータが不十分な状態で、SVM による知識利用を行うことは有効ではないと過去に示した（非特許文献 1 参照）。そこで、知識獲得手段 12 における学習が十分に進み、機械学習システム 1 の行動が安定し始めたタイミングで知識再構成手段 14 を動作させて獲得知識を再構成するため指標を設ける。

【0046】

まず、SVM による知識利用の前提として、サンプルデータが十分に多い、すなわち、知識獲得手段 12 においてパラメトリック表現されたクラス集合の生成に使用された学習済み入力の個数が十分に多い必要がある。したがって、学習済み入力の個数が所定数よりも多いことを指標に設定する。

【0047】

また、サンプルデータ数が十分に多くても、パラメトリック表現されたクラスの中に分散が大きいクラスが存在すると、SVM による十分な識別精度が得られないおそれがある。そこで、パラメトリック表現された各クラス（ルール）の分散が所定値よりも小さいことも指標に設定する。例えば、この指標を各ルールの構成要素である共分散行列 Σ により設定する。 σ の値は状態空間においてルールの範囲を表す。行動が収束し始めると σ の値が収束していく。よって、知識再構成手段 14 は、例えば、エピソード毎の各ルールの σ の平均 (σ_{eps}) を計算し、前エピソードの σ の平均 (σ_{eps-1}) との差が閾値以下の場合、SVM による知識再構成を行って、ノンパラメトリック表現されたクラス集合を生成する。そして、知識利用手段 14 は、当該ノンパラメトリック表現された

【0048】

【数 9】

$$\Sigma'_{eps} = \left| \Sigma_{eps} - \Sigma_{eps-1} \right| \quad (9)$$

【0049】

また、ルール判別可とする範囲を BR L の場合よりも広くするために新たな閾値 P_s ($P_s > P_{th}$) を設定する。閾値を広げることで新ルールの生成を抑制し、振る舞いの不安定化を防ぐことが期待される。以下に、機械学習システム 1 の行動選択の概要を示す。また、図 2 に、機械学習システム 1 による知識獲得および利用の概要を示す。

【0050】

・知覚した入力をもとに、BR L により知識探索を行うか SVM により知識利用を行うかを、 σ_{eps} の値と閾値との比較により決定する。

【0051】

・BR L により知識探索を行う場合、ランダム行動を出力し、罰を受けなければ新たなルールを作成する。

【0052】

・SVM により知識利用を行う場合、SVM により状態空間の再分割を行い、再分割し

10

20

30

40

50

た状態空間により判別されたルールの行動を出力する。

【0053】

なお、数式(9)は指標の一例であり本発明はこれに限定されない。例えば、 ρ の移動平均や加重平均を利用してもよい。

【0054】

多クラス分類SVMs

SVMは基本的には2クラスの識別問題を対象にして定式化されている。しかし、2クラスの判別モデルを組み合わせることで多クラス分類を可能にしている。組み合わせ方としてOne-versus-AllとOne-versus-Oneという2種類を取り上げる。One-versus-Allとは、全クラスに対して、ある一つのクラスとそれ以外のクラスに分ける識別平面を作成し、これらの識別平面のうち最も高い判別値を返すクラスを出力するという方法である。nクラスの問題の場合、識別平面の数はnとなる。一方、One-versus-Oneとは、各クラス毎に対応する識別平面を作成し、多数決により出力を決定する方法である。識別平面の数は $n(n-1)/2$ となる。知識利用手段14に用いるSVMに、この2種類の多クラス分類方式を導入する。

10

【0055】

SVMに用いるカーネル関数

知識再構成手段14で使用するSVMには高次元空間への写像により非線形分離を可能にするカーネルトリックを用いる。カーネルは以下に示す線形カーネル K_{line} 、多項式カーネル K_{poly} 、RBFカーネル K_{RBF} 、シグモイドカーネル K_{sig} を使用する。 u はサポートベクトル、 v は識別する特徴ベクトルを表す。

20

【0056】

【数10】

$$K_{line} = uv \quad (10)$$

【0057】

【数11】

$$K_{poly} = (\gamma uv + c)^d \quad (11)$$

30

【0058】

【数12】

$$K_{RBF} = \exp(-\gamma |u - v|^2) \quad (12)$$

【0059】

【数13】

$$K_{sig} = \tanh(\gamma uv + c) \quad (13)$$

40

【0060】

(計算機実験)

BRLによる知識獲得のタスクとして、二台のロボットによるピアノ運び問題を取り上げ、計算機実験を行う。ピアノ運び問題とは単体では搬送不可能な長尺物をロボットが協調しゴールまで搬送する問題である。狭い通路を通行するためには、二台のロボットが協

50

調し、フォーメーションを形成しなければならない。

【0061】

実験環境を図3(a)に示す。フィールドは四方を壁で囲まれており、初期位置からゴールラインまで移動するとタスク達成となる。ロボットは差動駆動型を用い、二輪の駆動輪を持つ。各ロボットは全方位カメラにより他のロボットの状態を知覚し、ロボット間で通信を行わない。ロボットの一度の意思決定を1ステップとし、タスクを達成するか400ステップ経過した時点でエピソードを更新し、ロボットを初期位置に戻す。学習成功の定義は、20エピソード連続でタスクを達成したときとする。500エピソード経過までを1試行とし、100試行繰り返す。なお、シミュレーション環境はオープンソース三次元物理エンジンODE (Open Dynamic Engine) により作成している。

10

【0062】

実験1

BRLのみを用いて学習を行う。学習が成功した試行には、図3(b)に示したような環境変化を行い、SVMによる知識利用の効果を検証するためSVMのみで行動選択を行う。環境変化は通路の幅が狭くなることであり、タスクの難易度が上昇するため、獲得知識を状態に応じて正確に判別することが求められる。ここで使用する多クラス分類SVMsは、2種類の多クラス分類方式と4種類のカーネル関数の組み合わせの8パターンを使用する。それぞれの組み合わせを表1で示すA~Hまでの記号で表す。学習が成功した試行中、環境変化のタスクを達成できた試行の成功率を観察する。

20

【0063】

【表1】

	1-All	1-1
$K_{\text{line}}[(C, \gamma)=(2^3, 2^{-2})]$	A	B
$K_{\text{poly}}[(C, \gamma, d, c)=(2^5, 2^{-3}, 2, 30)]$	C	D
$K_{\text{RBF}}[(C, \gamma)=(2^4, 2^0)]$	E	F
$K_{\text{sig}}[(C, \gamma, c)=(2^4, 2^{-4}, 0)]$	G	H

30

【0064】

実験2

実験1で成功率の高かった組み合わせパターンについて、機械学習システム1を用いて学習を行う。学習するのに要したエピソード数の推移を調べ、機械学習システム1の行動獲得に対する有効性を検証する。パラメータチューニングにより、SVMを使用する ϵ_{eps} の閾値を0.02と規定した。

【0065】

機械学習システム1の設定

入力は、 $I = \{r_0, \cos \theta_0, \sin \theta_0, r_1, \cos \theta_1, \sin \theta_1, r_2, \cos \theta_2, \sin \theta_2, r_3, \cos \theta_3, \sin \theta_3\}$ の11次元である。 r_i は対象物までの距離とその角度を、添字0, 1, 2は対象物がそれぞれゴールライン、最近傍の壁、第二近傍の壁、添字3は隣のロボットを示している。BRLの出力はロボットの左右のモータ回転速度 $O = \{m_0, m_1\}$ の2次元である。SVMの設定は、ライブラリLIBSVM (<http://www.csie.ntu.edu.tw/~cjlin/libsvm/index.html>参照) を改良したものを使用し、各パラメータはチューニングにより規定した。

40

【0066】

実験結果1

全100試行のうち、76試行において学習に成功した。この76試行に関して、A~

50

Hまでの各SVMによる行動選択をそれぞれ行い、76試行中タスクを達成できた試行の成功率を図4に示す。図4のBRLとはランダム行動をせず、学習成功時のみの知識をベイズ判別法により行動決定を行った場合のことである。C, D, E, Fの場合がタスク成功率が高く、それぞれ78%, 71%, 78%, 71%であった。BRLのみで環境変化のタスク達成率は51%であることから、カーネル関数が K_{poly} 、 K_{RBF} を用いた場合が判別性能が高いとわかる。BRLに比べ、SVMを使用したA~Hのすべての場合でタスク成功率が高いことからSVMの判別性能は比較的高いとも言える。また、One-versus-AllとOne-versus-Oneを K_{poly} と K_{RBF} の場合で比較すると、若干であるがOne-versus-Allの方が成功率が高い。これは一般的にOne-versus-Allはクラス数が1000などの多い場合に有効であるとも言われるが、クラス数(ルール数)が最大100のBRLにおいては有効性に差はない。

10

【0067】

実験結果2

実験1の結果から、実験2ではC, D, E, Fの4パターンの場合について、機械学習システム1を用いて学習を行った。100試行ずつ行った結果、C, D, E, Fはそれぞれ78, 79, 75, 79試行において学習に成功した。BRLのみの学習成功回数76試行と比較し、大きな差はないように思われる。学習成功までに要したエピソード数の平均と標準偏差を図5に示す。この結果にT検定を行ったところ、BRLのみと比較しDとFの場合に関しては有意水準1%において差があることが示された。つまり、DもしくはFの場合、学習するまでに要する収束速度が早いことが示された。

20

【0068】

実験1と2の結果から総合的に判断して、DもしくはFのパターンのSVMを用いることが行動獲得に有効であり、識別精度が高いと言える。

【産業上の利用可能性】

【0069】

本発明に係る機械学習システムおよび機械学習方法は、頑健性に優れ、環境変化にも対応可能であるため、マルチロボットシステムなどに有用である。また、ロボットに限らず、パターン認識における強化学習にも有用である。

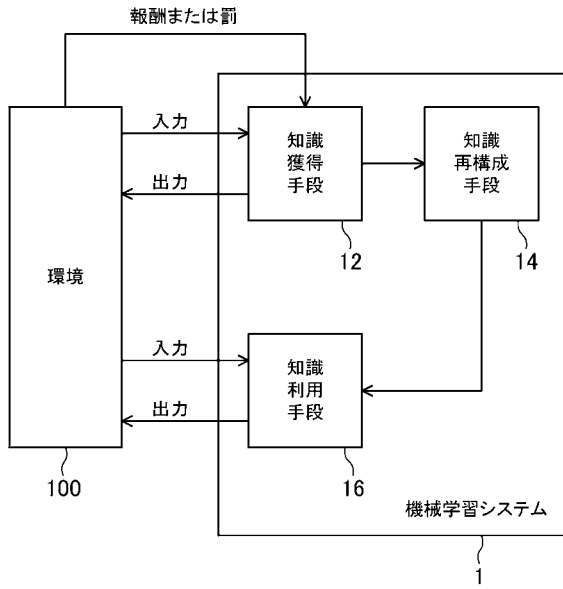
【符号の説明】

【0070】

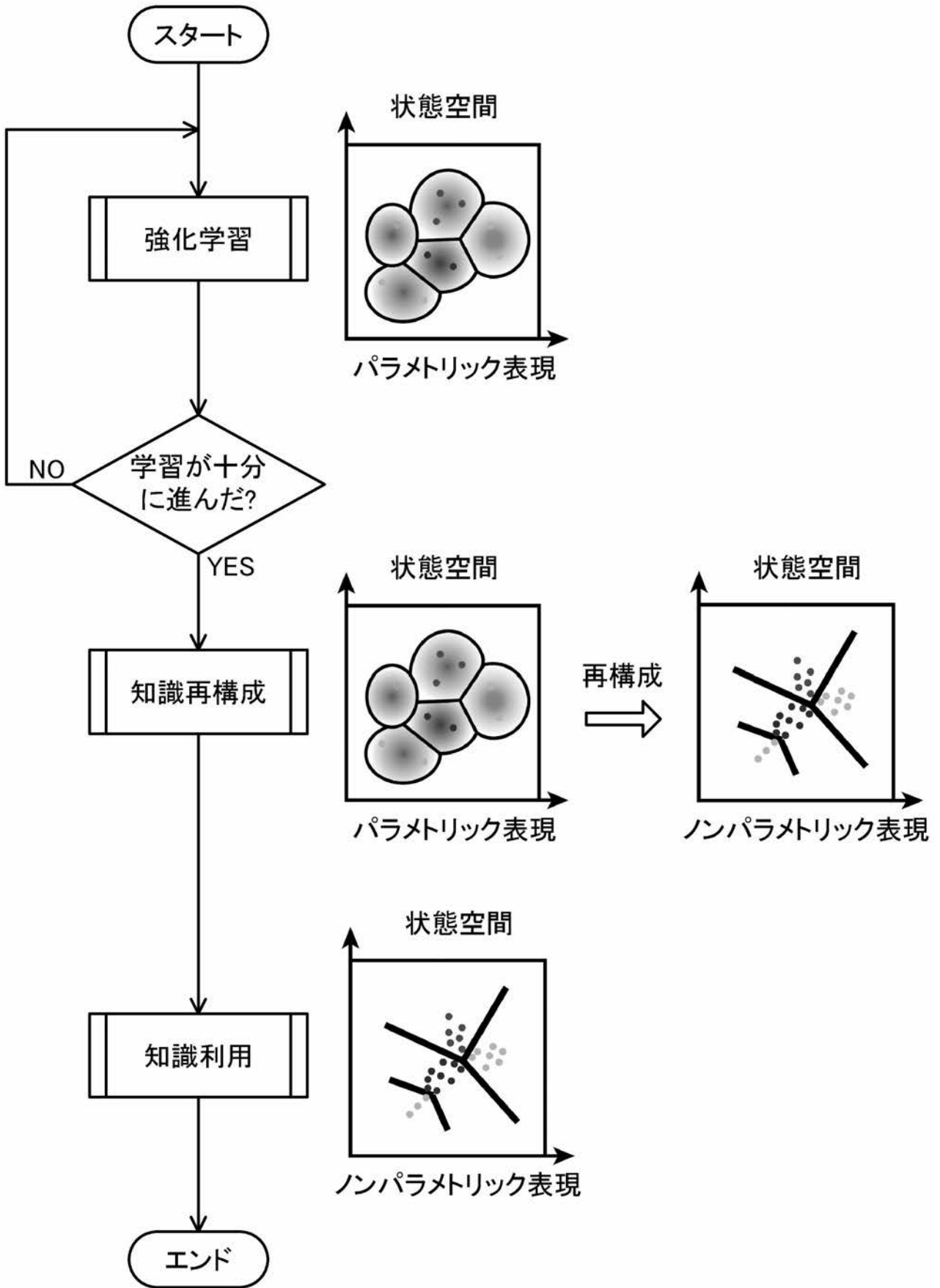
- 1 機械学習システム
- 12 知識獲得手段
- 14 知識再構成手段
- 16 知識利用手段
- 100 環境

30

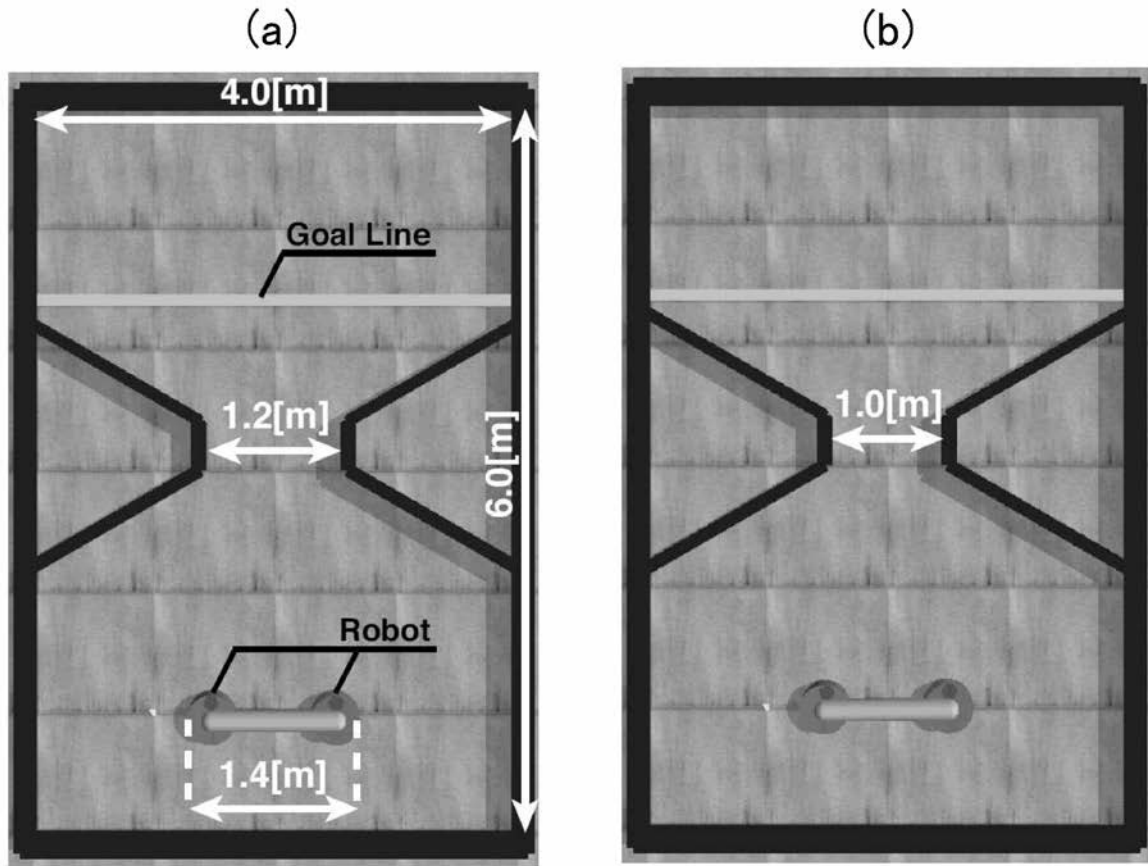
【 図 1 】



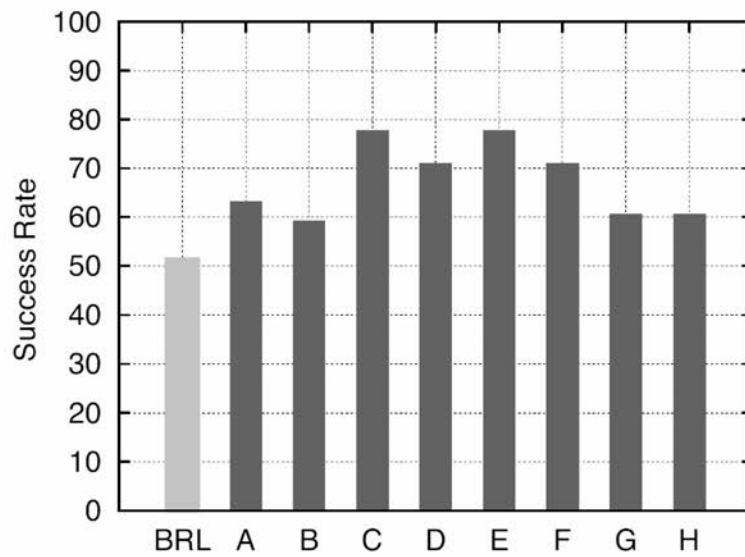
【図2】



【 図 3 】



【 図 4 】



【 図 5 】

