

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第5874292号
(P5874292)

(45) 発行日 平成28年3月2日(2016.3.2)

(24) 登録日 平成28年1月29日(2016.1.29)

(51) Int.Cl.		F I			
G06N 99/00	(2010.01)	G06N 99/00	150		
A63F 13/56	(2014.01)	A63F 13/56			
A63F 13/67	(2014.01)	A63F 13/67			

請求項の数 15 (全 60 頁)

(21) 出願番号	特願2011-224638 (P2011-224638)	(73) 特許権者	000002185
(22) 出願日	平成23年10月12日(2011.10.12)		ソニー株式会社
(65) 公開番号	特開2013-84175 (P2013-84175A)		東京都港区港南1丁目7番1号
(43) 公開日	平成25年5月9日(2013.5.9)	(74) 代理人	100095957
審査請求日	平成26年9月29日(2014.9.29)		弁理士 亀谷 美明
		(74) 代理人	100096389
			弁理士 金本 哲男
		(74) 代理人	100101557
			弁理士 萩原 康司
		(74) 代理人	100128587
			弁理士 松本 一騎
		(72) 発明者	小林 由幸
			東京都港区港南1丁目7番1号 ソニー株式会社内

最終頁に続く

(54) 【発明の名称】 情報処理装置、情報処理方法、及びプログラム

(57) 【特許請求の範囲】

【請求項1】

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成する報酬推定機生成部を備え、

前記報酬推定機生成部は、

複数の処理関数を組み合わせる複数の基底関数を生成する基底関数生成部と、

前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出部と、

前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出する推定関数算出部と、

を含み、

前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、

情報処理装置。

【請求項2】

エージェントがおかれた現在の状態を表す状態データと、当該エージェントが次にとりうる行動を表す行動データを前記報酬推定機に入力し、当該行動をとった結果として当該エージェントが得る報酬値を推定する報酬値推定部と、

前記現在の状態において前記エージェントが次にとりうる行動のうち、前記報酬値推定

部により推定された報酬値が最も高い値となる行動を選択する行動選択部と、
をさらに備える、

請求項 1 に記載の情報処理装置。

【請求項 3】

前記行動選択部による選択結果に基づいてエージェントを行動させる行動制御部と、
前記エージェントの行動に伴って更新される状態データ及び行動データを蓄積し、当該
行動の結果としてエージェントが得た報酬を表す報酬値と、蓄積した状態データ及び行動
データとを対応付けて前記行動履歴データに追加する履歴データ追加部と、
をさらに備える、

請求項 2 に記載の情報処理装置。

10

【請求項 4】

前記状態データ、前記行動データ、及び前記報酬値の組が前記行動履歴データに追加さ
れた場合、前記特徴量ベクトル算出部は、前記行動履歴データに含まれる全ての状態デー
タ及び行動データについて特徴量ベクトルを算出し、

前記情報処理装置は、特徴領空間において前記特徴量ベクトルにより示される座標点の
分布が所定の分布に近づくように前記行動履歴データに含まれる前記状態データ、前記行
動データ、及び前記報酬値の組を間引く分布調整部をさらに備える、

請求項 3 に記載の情報処理装置。

【請求項 5】

前記状態データ、前記行動データ、及び前記報酬値の組が前記行動履歴データに追加さ
れた場合、前記特徴量ベクトル算出部は、前記行動履歴データに含まれる全ての状態デー
タ及び行動データについて特徴量ベクトルを算出し、

前記情報処理装置は、特徴領空間において前記特徴量ベクトルにより示される座標点の
分布が所定の分布に近づくように前記行動履歴データに含まれる前記状態データ、前記行
動データ、及び前記報酬値の組のそれぞれに重みを設定する分布調整部をさらに備える、

請求項 3 に記載の情報処理装置。

20

【請求項 6】

前記分布調整部は、間引き後に残った前記状態データ、前記行動データ、及び前記報酬
値の組について、特徴領空間において前記特徴量ベクトルにより示される座標点の分布が
所定の分布に近づくように前記行動履歴データに含まれる前記状態データ、前記行動デー
タ、及び前記報酬値の組のそれぞれに重みを設定する、

請求項 4 に記載の情報処理装置。

30

【請求項 7】

前記行動履歴データを学習データとして用い、現在の時刻においてエージェントがおか
れた状態を表す状態データ及び現在の時刻においてエージェントがとる行動を表す行動デー
タから次の時刻におけるエージェントの状態を表す状態データを予測する予測機を機械
学習により生成する予測機生成部をさらに備え、

前記報酬値推定部は、

現在の時刻における状態データ及び行動データを前記予測機に入力して次の時刻におけ
るエージェントの状態を表す状態データを予測し、

前記次の時刻におけるエージェントの状態を表す状態データと、当該状態においてエー
ジェントがとりうる行動を表す行動データとを前記報酬推定機に入力して、当該行動をと
った結果として当該エージェントが得る報酬値を推定する、

請求項 2 ~ 6 のいずれか 1 項に記載の情報処理装置。

40

【請求項 8】

前記行動履歴データを学習データとして用い、現在の時刻においてエージェントがおか
れた状態を表す状態データ及び現在の時刻においてエージェントがとる行動を表す行動デー
タから次の時刻におけるエージェントの状態を表す状態データを予測する予測機を機械
学習により生成する予測機生成部をさらに備え、

前記報酬値推定部は、現在の時刻を時刻 t_0 とした場合に、

50

時刻 t_0 における状態データ及び行動データを前記予測機に入力して次の時刻 t_1 におけるエージェントの状態を表す状態データを予測する処理を実行し、

$k = 1 \sim n - 1$ ($n - 2$) について、時刻 t_k における状態データ及び時刻 t_k においてエージェントがとりうる行動を表す行動データを前記予測機に入力して時刻 t_{k+1} におけるエージェントの状態を表す状態データを予測する処理を逐次実行し、

予測した時刻 t_n におけるエージェントの状態を表す状態データと、当該状態においてエージェントがとりうる行動を表す行動データとを前記報酬推定機に入力して、当該行動をとった結果として当該エージェントが得る報酬値を推定する、

請求項 2 ~ 6 のいずれか 1 項に記載の情報処理装置。

【請求項 9】

前記報酬推定機生成部は、複数のエージェントの状態を表す状態データと、当該状態において各エージェントがとった行動を表す行動データと、当該行動の結果として各エージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成する、

請求項 1 ~ 8 のいずれか 1 項に記載の情報処理装置。

【請求項 10】

前記基底関数生成部は、遺伝的アルゴリズムに基づいて前記基底関数を更新し、

前記特徴量ベクトル算出部は、前記基底関数が更新された場合に、更新後の前記基底関数に前記状態データ及び前記行動データを入力して特徴量ベクトルを算出し、

前記推定関数算出部は、前記更新後の基底関数を用いて算出された特徴量ベクトルの入力に応じて前記報酬値を推定する推定関数を算出する、

請求項 1 ~ 9 のいずれか 1 項に記載の情報処理装置。

【請求項 11】

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかつた行動に高いスコアを与え、高い報酬を得たエージェントがとらなかつた行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するスコア算出部と、

前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するスコア推定機生成部と、

を備え、

前記スコア推定機生成部は、

複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成部と、

前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出部と、

前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰 / 判別学習により算出する推定関数算出部と、

を含み、

前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、

情報処理装置。

【請求項 12】

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成するステップを含み、

前記生成するステップは、

複数の処理関数を組み合わせて複数の基底関数を生成するステップと、

10

20

30

40

50

前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出するステップと、

前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出するステップと、
を含み、

前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、
情報処理方法。

【請求項13】

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するステップと、

前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するステップと、

を含み、

前記生成するステップは、

複数の処理関数を組み合わせて複数の基底関数を生成するステップと、

前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出するステップと、

前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習により算出するステップと、

を含み、

前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、
情報処理方法。

【請求項14】

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成する報酬推定機生成機能をコンピュータに実現させるためのプログラムであり、

前記報酬推定機生成機能は、

複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成機能と、

前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出機能と、

前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出する推定関数算出機能と、

を含み、

前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、
プログラム。

【請求項15】

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するスコア算出機能と、

前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態

10

20

30

40

50

データから行動毎のスコアを推定するスコア推定機を機械学習により生成するスコア推定機生成機能と、

をコンピュータに実現させるためのプログラムであり、

前記スコア推定機生成機能は、

複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成機能と、

前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出機能と、

前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習により算出する推定関数算出機能と、

を含み、

前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本技術は、情報処理装置、情報処理方法、及びプログラムに関する。

【背景技術】

【0002】

近年、定量的に特徴を決定づけることが難しい任意のデータ群から、そのデータ群の特徴量を機械的に抽出する手法に注目が集まっている。例えば、任意の音楽データを入力とし、その音楽データが属する音楽のジャンルを機械的に抽出するアルゴリズムを自動構築する手法が知られている。ジャズ、クラシック、ポップス等、音楽のジャンルは、楽器の種類や演奏形態に応じて定量的に決まるものではない。そのため、これまでは任意の音楽データが与えられたときに、その音楽データから機械的に音楽のジャンルを抽出することは一般的に難しいと考えられていた。

【0003】

しかし、実際には、音楽データに含まれる音程の組み合わせ、音程の組み合わせ方、楽器の種類や組み合わせ、メロディーラインやベースラインの構造等、様々な情報の組み合わせの中に、音楽のジャンルを分ける特徴が潜在的に含まれている。そのため、この特徴を抽出するアルゴリズム(以下、特徴量抽出機)を機械学習により自動構築できないか、という観点から特徴量抽出機の研究が行われた。その研究成果の一つとして、例えば、下記の特許文献1に記載された遺伝アルゴリズムに基づく特徴量抽出機の自動構築方法を挙げることができる。遺伝アルゴリズムとは、生物の進化過程に倣い、機械学習の過程で、選択、交差、突然変異の要素を考慮したものを言う。

【0004】

同文献に記載の特徴量抽出機自動構築アルゴリズムを利用することにより、任意の音楽データから、その音楽データが属する音楽のジャンルを抽出する特徴量抽出機を自動構築することができるようになる。また、同文献に記載の特徴量抽出機自動構築アルゴリズムは、非常に汎用性が高く、音楽データに限らず、任意のデータ群から、そのデータ群の特徴量を抽出する特徴量抽出機を自動構築することができる。そのため、同文献に記載の特徴量抽出機自動構築アルゴリズムは、音楽データや映像データのような人工的なデータの特徴量解析、自然界に存在する様々な観測の特徴量解析への応用が期待されている。

【先行技術文献】

【特許文献】

【0005】

【特許文献1】特開2009-48266号公報

【発明の概要】

【発明が解決しようとする課題】

【0006】

ところで、本件発明者は、同文献に記載の技術に対して更なる工夫を施すことで、エー

10

20

30

40

50

ジェントを賢く行動させるアルゴリズムを自動構築する技術に発展させられないか検討を行ってきた。その検討の中で、本件発明者は、ある状態におかれたエージェントがとりうる行動の中で選択すべき行動を決定するための思考ルーチンを自動構築する技術に注目した。本技術は、このような技術に関するものであり、エージェントがとるべき行動を選択する際に決め手となる情報を出力する推定機を自動構築することが可能な、新規かつ改良された情報処理装置、情報処理方法、及びプログラムを提供することを意図している。

【課題を解決するための手段】

【0007】

本技術のある観点によれば、エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成する報酬推定機生成部を備え、前記報酬推定機生成部は、複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成部と、前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出部と、前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出する推定関数算出部と、を含み、前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、情報処理装置が提供される。

10

【0008】

また、本技術の別の観点によれば、エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するスコア算出部と、前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するスコア推定機生成部と、を備え、前記スコア推定機生成部は、複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成部と、前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出部と、前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習により算出する推定関数算出部と、を含み、前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、情報処理装置が提供される。

20

30

【0009】

また、本技術の別の観点によれば、エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成するステップを含み、前記生成するステップは、複数の処理関数を組み合わせて複数の基底関数を生成するステップと、前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出するステップと、前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出するステップと、を含み、前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、情報処理方法が提供される。

40

【0010】

また、本技術の別の観点によれば、エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントが

50

とった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するステップと、前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するステップと、を含み、前記生成するステップは、複数の処理関数を組み合わせて複数の基底関数を生成するステップと、前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出するステップと、前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習により算出するステップと、を含み、前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、情報処理方法が提供される。

【0011】

また、本技術の別の観点によれば、エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成する報酬推定機生成機能をコンピュータに実現させるためのプログラムであり、前記報酬推定機生成機能は、複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成機能と、前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出機能と、前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出する推定関数算出機能と、を含み、前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、プログラムが提供される。

【0012】

また、本技術の別の観点によれば、エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するスコア算出機能と、前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するスコア推定機生成機能と、をコンピュータに実現させるためのプログラムであり、前記スコア推定機生成機能は、複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成機能と、前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出機能と、前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習により算出する推定関数算出機能と、を含み、前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、プログラムが提供される。

【0013】

また、本技術の別の観点によれば、上記のプログラムが記録された、コンピュータにより読み取り可能な記録媒体が提供される。

【発明の効果】

【0014】

以上説明したように本技術によれば、エージェントがとるべき行動を選択する際に決め手となる情報を出力する推定機を自動構築することが可能になる。

【図面の簡単な説明】

【0015】

【図1】推定機の自動構築方法について説明するための説明図である。

【図2】推定機の自動構築方法について説明するための説明図である。

【図3】推定機の自動構築方法について説明するための説明図である。

【図4】推定機の自動構築方法について説明するための説明図である。

10

20

30

40

50

- 【図5】推定機の自動構築方法について説明するための説明図である。
- 【図6】推定機の自動構築方法について説明するための説明図である。
- 【図7】推定機の自動構築方法について説明するための説明図である。
- 【図8】推定機の自動構築方法について説明するための説明図である。
- 【図9】推定機の自動構築方法について説明するための説明図である。
- 【図10】推定機の自動構築方法について説明するための説明図である。
- 【図11】推定機の自動構築方法について説明するための説明図である。
- 【図12】推定機の自動構築方法について説明するための説明図である。
- 【図13】オンライン学習に基づく推定機の自動構築方法について説明するための説明図である。 10
- 【図14】データセットの統合方法について説明するための説明図である。
- 【図15】データセットの統合方法について説明するための説明図である。
- 【図16】データセットの統合方法について説明するための説明図である。
- 【図17】データセットの統合方法について説明するための説明図である。
- 【図18】データセットの統合方法について説明するための説明図である。
- 【図19】データセットの統合方法について説明するための説明図である。
- 【図20】データセットの統合方法について説明するための説明図である。
- 【図21】データセットの統合方法について説明するための説明図である。
- 【図22】データセットの統合方法について説明するための説明図である。
- 【図23】データセットの統合方法について説明するための説明図である。 20
- 【図24】データセットの統合方法について説明するための説明図である。
- 【図25】データセットの統合方法について説明するための説明図である。
- 【図26】データセットの統合方法について説明するための説明図である。
- 【図27】データセットの統合方法について説明するための説明図である。
- 【図28】データセットの統合方法について説明するための説明図である。
- 【図29】データセットの統合方法について説明するための説明図である。
- 【図30】データセットの統合方法について説明するための説明図である。
- 【図31】データセットの統合方法について説明するための説明図である。
- 【図32】データセットの統合方法について説明するための説明図である。
- 【図33】データセットの統合方法について説明するための説明図である。 30
- 【図34】思考ルーチンの構成について説明するための説明図である。
- 【図35】思考ルーチンの構成について説明するための説明図である。
- 【図36】思考ルーチンの構成について説明するための説明図である。
- 【図37】思考ルーチンの構成について説明するための説明図である。
- 【図38】思考ルーチンの構築方法について説明するための説明図である。
- 【図39】情報処理装置10の機能構成例について説明するための説明図である。
- 【図40】情報処理装置10の機能構成例について説明するための説明図である。
- 【図41】効率的な報酬推定機の構築方法について説明するための説明図である。
- 【図42】効率的な報酬推定機の構築方法について説明するための説明図である。
- 【図43】アクションスコア推定機を用いた思考ルーチンの構成について説明するための説明図である。 40
- 【図44】アクションスコア推定機を用いた思考ルーチンの構成について説明するための説明図である。
- 【図45】予測機を用いた報酬の推定方法について説明するための説明図である。
- 【図46】予測機を用いた報酬の推定方法について説明するための説明図である。
- 【図47】予測機を用いた報酬の推定方法について説明するための説明図である。
- 【図48】「三目並べ」への応用について説明するための説明図である。
- 【図49】「三目並べ」への応用について説明するための説明図である。
- 【図50】「三目並べ」への応用について説明するための説明図である。
- 【図51】「三目並べ」への応用について説明するための説明図である。 50

- 【図52】「三目並べ」への応用について説明するための説明図である。
- 【図53】「三目並べ」への応用について説明するための説明図である。
- 【図54】「三目並べ」への応用について説明するための説明図である。
- 【図55】「対戦ゲーム」への応用について説明するための説明図である。
- 【図56】「対戦ゲーム」への応用について説明するための説明図である。
- 【図57】「対戦ゲーム」への応用について説明するための説明図である。
- 【図58】「対戦ゲーム」への応用について説明するための説明図である。
- 【図59】「対戦ゲーム」への応用について説明するための説明図である。
- 【図60】「対戦ゲーム」への応用について説明するための説明図である。
- 【図61】「対戦ゲーム」への応用について説明するための説明図である。 10
- 【図62】「対戦ゲーム」への応用について説明するための説明図である。
- 【図63】「五目並べ」への応用について説明するための説明図である。
- 【図64】「五目並べ」への応用について説明するための説明図である。
- 【図65】「ポーカー」への応用について説明するための説明図である。
- 【図66】「ポーカー」への応用について説明するための説明図である。
- 【図67】「ポーカー」への応用について説明するための説明図である。
- 【図68】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図69】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図70】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図71】「ロールプレイングゲーム」への応用について説明するための説明図である。 20
- 【図72】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図73】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図74】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図75】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図76】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図77】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図78】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図79】「ロールプレイングゲーム」への応用について説明するための説明図である。
- 【図80】情報処理装置の機能を実現することが可能なハードウェア構成例について説明するための説明図である。 30
- 【発明を実施するための形態】
- 【0016】
- 以下に添付図面を参照しながら、本技術に係る好適な実施の形態について詳細に説明する。なお、本明細書及び図面において、実質的に同一の機能構成を有する構成要素については、同一の符号を付することにより重複説明を省略する。
- 【0017】
- [説明の流れについて]
- ここで、以下に記載する説明の流れについて簡単に述べる。
- 【0018】
- まず、本実施形態に係る基盤技術について説明する。具体的には、まず、図1～図12を参照しながら、推定機の自動構築方法について説明する。次いで、図13を参照しながら、オンライン学習に基づく推定機の自動構築方法について説明する。 40
- 【0019】
- 次いで、図14～図16を参照しながら、データセットの統合方法について説明する。次いで、図17～図23を参照しながら、効率的なデータセットのサンプリング方法について説明する。次いで、図24～図27を参照しながら、効率的な重み付け方法について説明する。次いで、図28を参照しながら、効率的なデータセットのサンプリング方法及び重み付け方法を組み合わせる方法について説明する。次いで、図29～図33を参照しながら、その他のデータセットのサンプリング方法及び重み付け方法について説明する。
- 【0020】 50

次いで、図34～図38を参照しながら、思考ルーチンの構成及び思考ルーチンの構築方法について説明する。次いで、図39及び図40を参照しながら、本実施形態に係る情報処理装置10の機能構成について説明する。次いで、図41及び図42を参照しながら、効率的な報酬推定機の構築方法について説明する。次いで、図43及び図44を参照しながら、アクションスコア推定機を用いた思考ルーチンの構成について説明する。次いで、図45～図47を参照しながら、予測機を用いた報酬の推定方法について説明する。

【0021】

次いで、図48～図54を参照しながら、本実施形態に係る技術を「三目並べ」へ応用する方法について説明する。次いで、図55～図62を参照しながら、本実施形態に係る技術を「対戦ゲーム」へ応用する方法について説明する。次いで、図63及び図64を参照しながら、本実施形態に係る技術を「五目並べ」へ応用する方法について説明する。次いで、図65～図67を参照しながら、本実施形態に係る技術を「ポーカー」へ応用する方法について説明する。次いで、図68～図79を参照しながら、本実施形態に係る技術を「ロールプレイングゲーム」へ応用する方法について説明する。

【0022】

次いで、図80を参照しながら、本実施形態に係る情報処理装置10の機能を実現することが可能なハードウェア構成例について説明する。最後に、同実施形態の技術的思想について纏め、当該技術的思想から得られる作用効果について簡単に説明する。

【0023】

(説明項目)

1：基盤技術

- 1 - 1：推定機の自動構築方法
 - 1 - 1 - 1：推定機の構成
 - 1 - 1 - 2：構築処理の流れ
- 1 - 2：オンライン学習について
- 1 - 3：学習用データの統合方法
 - 1 - 3 - 1：特徴量空間における学習用データの分布と推定機の精度
 - 1 - 3 - 2：データ統合時にサンプリングする構成
 - 1 - 3 - 3：データ統合時に重み付けする構成
 - 1 - 3 - 4：データ統合時にサンプリング及び重み付けする構成
- 1 - 4：効率的なサンプリング/重み付け方法
 - 1 - 4 - 1：サンプリング方法
 - 1 - 4 - 2：重み付け方法
 - 1 - 4 - 3：組み合わせ方法
- 1 - 5：サンプリング処理及び重み付け処理に関する変形例
 - 1 - 5 - 1：変形例1（距離に基づく処理）
 - 1 - 5 - 2：変形例2（クラスタリングに基づく処理）
 - 1 - 5 - 3：変形例3（密度推定手法に基づく処理）

2：実施形態

- 2 - 1：思考ルーチンの自動構築方法
 - 2 - 1 - 1：思考ルーチンとは
 - 2 - 1 - 2：思考ルーチンの構成
 - 2 - 1 - 3：報酬推定機の構築方法
- 2 - 2：情報処理装置10の構成
- 2 - 3：効率的な推定報酬機の構築方法
- 2 - 4：（変形例1）アクションスコア推定機を用いる思考ルーチン
- 2 - 5：（変形例2）予測機を用いた報酬の推定
 - 2 - 5 - 1：予測機の構築方法
 - 2 - 5 - 2：報酬の推定方法
- 2 - 6：（変形例3）複数エージェントの同時学習

10

20

30

40

50

3 : 応用例

- 3 - 1 : 「三目並べ」への応用
- 3 - 2 : 「対戦ゲーム」への応用
- 3 - 3 : 「五目並べ」への応用
- 3 - 4 : 「ポーカー」への応用
- 3 - 5 : 「ロールプレイングゲーム」への応用

4 : ハードウェア構成例

5 : まとめ

【 0 0 2 4 】

< 1 : 基盤技術 >

後述する実施形態は、推定機の自動構築方法に関する。また、同実施形態は、推定機の構築に用いる学習用データを追加できるようにする仕組み（以下、オンライン学習）に関する。そこで、同実施形態に係る技術について詳細に説明するに先立ち、推定機の自動構築方法及びオンライン学習方法（以下、基盤技術）について説明する。なお、以下では遺伝アルゴリズムに基づく推定機の自動構築方法を例に挙げて説明を進めるが、同実施形態に係る技術の適用範囲はこれに限定されない。

【 0 0 2 5 】

[1 - 1 : 推定機の自動構築方法]

推定機の自動構築方法について説明する。

【 0 0 2 6 】

(1 - 1 - 1 : 推定機の構成)

はじめに、図 1 ~ 図 3 を参照しながら、推定機の構成について説明する。図 1 は、推定機を利用するシステムのシステム構成例を示した説明図である。また、図 2 は、推定機の構築に利用する学習用データの構成例を示した説明図である。そして、図 3 は、推定機の構造及び構築方法の概要を示した説明図である。

【 0 0 2 7 】

まず、図 1 を参照する。図 1 に示すように、推定機の構築及び推定値の算出は、例えば、情報処理装置 1 0 により実行される。情報処理装置 1 0 は、学習用データ (X_1, t_1) , ..., (X_N, t_N) を利用して推定機を構築する。また、情報処理装置 1 0 は、構築した推定機を利用して入力データ X から推定値 y を算出する。この推定値 y は、入力データ X の認識に利用される。例えば、推定値 y が所定の閾値 T_h より大きい場合に認識結果 YES が得られ、推定値 y が所定の閾値 T_h より小さい場合に認識結果 NO が得られる。

【 0 0 2 8 】

図 2 を参照しながら、より具体的に推定機の構成について考えてみよう。図 2 に例示した学習用データの集合は、“海”の画像を認識する画像認識機の構築に利用されるものである。この場合、情報処理装置 1 0 により構築される推定機は、入力された画像の“海らしさ”を表す推定値 y を出力するものとなる。図 2 に示すように、学習用データは、データ X_k と目的変数 t_k とのペア（但し、 $k = 1 \sim N$ ）により構成される。データ X_k は、 k 番目の画像データ（画像 # k ）である。また、目的変数 t_k は、画像 # k が“海”の画像である場合に 1、画像 # k が“海”の画像でない場合に 0 となる変数である。

【 0 0 2 9 】

図 2 の例では、画像 # 1 が“海”の画像であり、画像 # 2 が“海”の画像であり、...、画像 # N が“海”の画像でない。この場合、 $t_1 = 1$ 、 $t_2 = 1$ 、...、 $t_N = 0$ となる。この学習用データが入力されると、情報処理装置 1 0 は、入力された学習用データに基づく機械学習により、入力された画像の“海らしさ”を表す推定値 y を出力する推定機を構築する。この推定値 y は、入力された画像の“海らしさ”が高いほど 1 に近づき、“海らしさ”が低いほど 0 に近づく値である。

【 0 0 3 0 】

また、新たに入力データ X （画像 X ）が入力されると、情報処理装置 1 0 は、学習用データの集合を利用して構築された推定機に画像 X を入力し、画像 X の“海らしさ”を表す

10

20

30

40

50

推定値 y を算出する。この推定値 y を利用すると、画像 X が“海”の画像であるか否かを認識することが可能になる。例えば、推定値 y が所定の閾値 T_h の場合、入力された画像 X が“海”の画像であると認識される。一方、推定値 $y < T_h$ の場合、入力された画像 X が“海”の画像でないと認識される。

【0031】

本実施形態は、上記のような推定機を自動構築する技術に関する。なお、ここでは画像認識機の構築に利用される推定機について説明したが、本実施形態に係る技術は、様々な推定機の自動構築方法に適用することができる。例えば、言語解析機の構築に適用することもできるし、楽曲のメロディーラインやコード進行などを解析する音楽解析機の構築にも適用することができる。さらに、蝶の動きや雲の流れなどの自然現象を再現したり、自然の振る舞いを予測したりする動き予測機の構築などにも適用することができる。

10

【0032】

例えば、特開2009-48266号公報、特願2010-159598号明細書、特願2010-159597号明細書、特願2009-277083号明細書、特願2009-277084号明細書などに記載のアルゴリズムに適用することができる。また、AdaBoostなどのアンサンブル学習手法や、SVMやSVRなどのカーネルを用いた学習手法などにも適用できる。AdaBoostなどのアンサンブル学習手法に適用する場合、弱学習機(Weak Learner)が後述する基底関数に対応する。また、SVMやSVRなどの学習手法に適用する場合、カーネルが後述する基底関数に対応する。なお、SVMはSupport Vector Machine、SVRはSupport Vector Regression、RVMはRelevance Vector Machineの略である。

20

【0033】

ここで、図3を参照しながら、推定機の構造について説明する。図3に示すように、推定機は、基底関数リスト (f_1, \dots, f_M) 及び推定関数 f により構成される。基底関数リスト (f_1, \dots, f_M) は、 M 個の基底関数 f_k ($k = 1 \sim M$) を含む。また、基底関数 f_k は、入力データ X の入力に応じて特徴量 z_k を出力する関数である。さらに、推定関数 f は、 M 個の特徴量 z_k ($k = 1 \sim M$) を要素として含む特徴量ベクトル $Z = (z_1, \dots, z_M)$ の入力に応じて推定値 y を出力する関数である。基底関数 f_k は、予め用意された1又は複数の処理関数を組み合わせて生成される。

30

【0034】

処理関数としては、例えば、三角関数、指数関数、四則演算、デジタルフィルタ、微分演算、中央値フィルタ、正規化演算、ホワイトノイズの付加処理、画像処理フィルタなどが利用可能である。例えば、入力データ X が画像の場合、ホワイトノイズの付加処理 $AddWhiteNoise()$ 、中央値フィルタ $Median()$ 、ぼかし処理 $Blur()$ を組み合わせた基底関数 $f_j(X) = AddWhiteNoise(Median(Blur(X)))$ などが利用される。この基底関数 f_j は、入力データ X に対し、ぼかし処理、中央値フィルタ処理、及びホワイトノイズの付加処理を順次施すことを意味する。

【0035】

(1-1-2: 構築処理の流れ)

40

さて、基底関数 f_k ($k = 1 \sim M$) の構成、基底関数リストの構成、推定関数 f の構成は、学習用データに基づく機械学習により決定される。以下、この機械学習による推定機の構築処理について、より詳細に説明する。

【0036】

(全体構成)

まず、図4を参照しながら、全体的な処理の流れについて説明する。図4は、全体的な処理の流れについて説明するための説明図である。なお、以下で説明する処理は、情報処理装置10により実行される。

【0037】

図4に示すように、まず、情報処理装置10に学習用データが入力される(S101)

50

。なお、学習用データとしては、データ X と目的変数 t の組が入力される。学習用データが入力されると、情報処理装置10は、処理関数を組み合わせて基底関数を生成する(S102)。次いで、情報処理装置10は、基底関数にデータ X を入力して特徴量ベクトル Z を算出する(S103)。次いで、情報処理装置10は、基底関数の評価及び推定関数の生成を行う(S104)。

【0038】

次いで、情報処理装置10は、所定の終了条件を満たしたか否かを判定する(S105)。所定の終了条件を満たした場合、情報処理装置10は、処理をステップS106に進める。一方、所定の終了条件を満たしていない場合、情報処理装置10は、処理をステップS102に戻し、ステップS102～S104の処理を繰り返し実行する。処理をステップS106に進めた場合、情報処理装置10は、推定関数を出力する(S106)。上記の通り、ステップS102～S104の処理は、繰り返し実行される。そこで、以下の説明においては、第 回目の繰り返し処理においてステップS102で生成される基底関数を第 世代の基底関数と呼ぶことにする。

10

【0039】

(基底関数の生成(S102))

ここで、図5～図10を参照しながら、ステップS102の処理(基底関数の生成)について、より詳細に説明する。

【0040】

まず、図5を参照する。図5に示すように、情報処理装置10は、現在の世代が2世代目以降であるか否かを判定する(S111)。つまり、情報処理装置10は、現在実行しようとしているステップS102の処理が第2回目以降の繰り返し処理であるか否かを判定する。2世代目以降である場合、情報処理装置10は、処理をステップS113に進める。一方、2世代目以降でない場合(第1世代である場合)、情報処理装置10は、処理をステップS112に進める。処理をステップS112に進めた場合、情報処理装置10は、基底関数をランダムに生成する(S112)。一方、処理をステップS113に進めた場合、情報処理装置10は、基底関数を進化的に生成する(S113)。そして、情報処理装置10は、ステップS112又はS113の処理が完了すると、ステップS102の処理を終了する。

20

【0041】

(S112:基底関数をランダムに生成)

次に、図6及び図7を参照しながら、ステップS112の処理について、より詳細に説明する。ステップS112の処理は、第1世代の基底関数を生成する処理に関する。

30

【0042】

まず、図6を参照する。図6に示すように、情報処理装置10は、基底関数のインデックス m ($m=0\sim M-1$)に関する処理ループを開始する(S121)。次いで、情報処理装置10は、基底関数 $\phi_m(x)$ をランダムに生成する(S122)。次いで、情報処理装置10は、基底関数のインデックス m が $M-1$ に達したか否かを判定し、基底関数のインデックス m が $M-1$ に達していない場合、情報処理装置10は、基底関数のインデックス m をインクリメントしてステップS121に処理を戻す(S124)。一方、基底関数のインデックス m が $m=M-1$ の場合、情報処理装置10は、処理ループを終了する(S124)。ステップS124で処理ループを終了すると、情報処理装置10は、ステップS112の処理を完了する。

40

【0043】

(ステップS122の詳細)

次に、図7を参照しながら、ステップS122の処理について、より詳細に説明する。

【0044】

ステップS122の処理を開始すると、図7に示すように、情報処理装置10は、基底関数のプロトタイプをランダムに決定する(S131)。プロトタイプとしては、既に例示した処理関数の他、線形項、ガウシアンカーネル、シグモイドカーネルなどの処理関数

50

が利用可能である。次いで、情報処理装置 10 は、決定したプロトタイプのパラメータをランダムに決定し、基底関数を生成する (S 1 3 2)。

【0045】

(S 1 1 3 : 基底関数を進化的に生成)

次に、図 8 ~ 図 10 を参照しながら、ステップ S 1 1 3 の処理について、より詳細に説明する。ステップ S 1 1 3 の処理は、第 t 世代 (t) の基底関数を生成する処理に関する。従って、ステップ S 1 1 3 を実行する際には、第 $t-1$ 世代の基底関数 $m, t-1$ ($m = 1 \sim M$) 及び当該基底関数 $m, t-1$ の評価値 $v_{m, t-1}$ が得られている。

【0046】

まず、図 8 を参照する。図 8 に示すように、情報処理装置 10 は、基底関数の数 M を更新する (S 1 4 1)。つまり、情報処理装置 10 は、第 t 世代の基底関数の数 M を決定する。次いで、情報処理装置 10 は、第 $t-1$ 世代の基底関数 $m, t-1$ ($m = 1 \sim M$) に対する評価値 $v_{t-1} = \{ v_{1, t-1}, \dots, v_{M, t-1} \}$ に基づき、第 $t-1$ 世代の基底関数の中から e 個の有用な基底関数を選択して第 t 世代の基底関数 $1, t, \dots, e, t$ に設定する (S 1 4 2)。

10

【0047】

次いで、情報処理装置 10 は、残り ($M - e$) 個の基底関数 $e+1, t, \dots, M, t$ を生成する方法を交差、突然変異、ランダム生成の中からランダムに選択する (S 1 4 3)。交差を選択した場合、情報処理装置 10 は、処理をステップ S 1 4 4 に進める。また、突然変異を選択した場合、情報処理装置 10 は、処理をステップ S 1 4 5 に進める。そして、ランダム生成を選択した場合、情報処理装置 10 は、処理をステップ S 1 4 6 に進める。

20

【0048】

処理をステップ S 1 4 4 に進めた場合、情報処理装置 10 は、ステップ S 1 4 2 で選択された基底関数 $1, t, \dots, e, t$ の中から選択された基底関数を交差させて新たな基底関数 m', t ($m' = e+1$) を生成する (S 1 4 4)。また、処理をステップ S 1 4 5 に進めた場合、情報処理装置 10 は、ステップ S 1 4 2 で選択された基底関数 $1, t, \dots, e, t$ の中から選択された基底関数を突然変異させて新たな基底関数 m', t ($m' = e+1$) を生成する (S 1 4 5)。一方、処理をステップ S 1 4 6 に進めた場合、情報処理装置 10 は、ランダムに新たな基底関数 m', t ($m' = e+1$) を生成する (S 1 4 6)。

30

【0049】

ステップ S 1 4 4、S 1 4 5、S 1 4 6 のいずれかの処理を終えると、情報処理装置 10 は、処理をステップ S 1 4 7 に進める。処理をステップ S 1 4 7 に進めると、情報処理装置 10 は、第 t 世代の基底関数が M 個 ($M = M$) に達したか否かを判定する (S 1 4 7)。第 t 世代の基底関数が M 個に達していない場合、情報処理装置 10 は、処理を再びステップ S 1 4 3 に戻す。一方、第 t 世代の基底関数が M 個に達した場合、情報処理装置 10 は、ステップ S 1 1 3 の処理を終了する。

【0050】

(S 1 4 4 の詳細 : 交差)

40

次に、図 9 を参照しながら、ステップ S 1 4 4 の処理について、より詳細に説明する。

【0051】

ステップ S 1 4 4 の処理を開始すると、図 9 に示すように、情報処理装置 10 は、ステップ S 1 4 2 で選択された基底関数 $1, t, \dots, e, t$ の中から同じプロトタイプを持つ基底関数をランダムに 2 つ選択する (S 1 5 1)。次いで、情報処理装置 10 は、選択した 2 つの基底関数が持つパラメータを交差させて新たな基底関数を生成する (S 1 5 2)。

【0052】

(S 1 4 5 の詳細 : 突然変異)

次に、図 10 を参照しながら、ステップ S 1 4 5 の処理について、より詳細に説明する

50

【0053】

ステップS145の処理を開始すると、図10に示すように、情報処理装置10は、ステップS142で選択された基底関数 ϕ_1, \dots, ϕ_e の中から基底関数をランダムに1つ選択する(S161)。次いで、情報処理装置10は、選択した基底関数を持つパラメータの一部をランダムに変更して新たな基底関数を生成する(S162)。

【0054】

(S146の詳細：ランダム生成)

次に、図7を参照しながら、ステップS146の処理について、より詳細に説明する。

【0055】

ステップS122の処理を開始すると、図7に示すように、情報処理装置10は、基底関数のプロトタイプをランダムに決定する(S131)。プロトタイプとしては、既に例示した処理関数の他、線形項、ガウシアンカーネル、シグモイドカーネルなどの処理関数が利用可能である。次いで、情報処理装置10は、決定したプロトタイプのパラメータをランダムに決定し、基底関数を生成する(S132)。

【0056】

以上、ステップS102の処理(基底関数の生成)について、より詳細に説明した。

【0057】

(基底関数の計算(S103))

次に、図11を参照しながら、ステップS103の処理(基底関数の計算)について、より詳細に説明する。

【0058】

図11に示すように、情報処理装置10は、学習用データに含まれる*i*番目のデータ $X^{(i)}$ のインデックス*i*に関する処理ループを開始する(S171)。例えば、学習用データとして*N*個のデータの組 $\{X^{(1)}, \dots, X^{(N)}\}$ が入力された場合には、 $i = 1 \sim N$ に関して処理ループが実行される。次いで、情報処理装置10は、基底関数 ϕ_m のインデックス*m*に関する処理ループを開始する(S172)。例えば、*M*個の基底関数を生成した場合には、 $m = 1 \sim M$ に関して処理ループが実行される。

【0059】

次いで、情報処理装置10は、特徴量 $z_{mi} = \phi_m(X^{(i)})$ を計算する(S173)。次いで、情報処理装置10は、処理をステップS174に進め、基底関数のインデックス*m*に関する処理ループを続ける。そして、情報処理装置10は、基底関数のインデックス*m*に関する処理ループが終了すると、処理をステップS175に進め、インデックス*i*に関する処理ループを続ける。インデックス*i*に関する処理ループが終了した場合、情報処理装置10は、ステップS103の処理を終了する。

【0060】

以上、ステップS103の処理(基底関数の計算)について、より詳細に説明した。

【0061】

(基底関数の評価・推定関数の生成(S104))

次に、図12を参照しながら、ステップS104の処理(基底関数の評価・推定関数の生成)について、より詳細に説明する。

【0062】

図12に示すように、情報処理装置10は、AIC基準の増減法に基づく回帰/判別学習により推定関数のパラメータ $w = \{w_0, \dots, w_M\}$ を算出する(S181)。つまり、情報処理装置10は、特徴量 $z_{mi} = \phi_m(X^{(i)})$ と目的変数 $t^{(i)}$ の組 ($i = 1 \sim N$) が推定関数 f によりフィッティングされるように、回帰/判別学習によりベクトル $w = \{w_0, \dots, w_M\}$ を求める。但し、推定関数 $f(x)$ は、 $f(x) = \sum_m w_m \phi_m(x) + w_0$ であるとする。次いで、情報処理装置10は、パラメータ w が0となる基底関数の評価値 v を0に設定し、それ以外の基底関数の評価値 v を1に設定する(S182)。つまり、評価値 v が1の基底関数は有用な基底関数である。

10

20

30

40

50

【 0 0 6 3 】

以上、ステップ S 1 0 4 の処理（基底関数の評価・推定関数の生成）について、より詳細に説明した。

【 0 0 6 4 】

推定機の構築に係る処理の流れは上記の通りである。このように、ステップ S 1 0 2 ~ S 1 0 4 の処理が繰り返し実行され、基底関数が進化的手法により逐次更新されることにより推定精度の高い推定関数が得られる。つまり、上記の方法を適用することで、高性能な推定機を自動構築することができる。

【 0 0 6 5 】

[1 - 2 : オンライン学習について]

さて、上記のように、機械学習により推定機を自動構築するアルゴリズムの場合、学習用データの数が多いほど、構築される推定機の性能が高くなる。そのため、可能な限り多くの学習用データを利用して推定機を構築するのが好ましい。また、後述する実施形態に係る技術においては、学習用データを追加する仕組みが利用される。そこで、学習用データを追加できるようにする新たな仕組み（以下、オンライン学習）について紹介する。

【 0 0 6 6 】

オンライン学習に係る推定機の構築は、図 1 3 に示すような処理の流れに沿って行われる。図 1 3 に示すように、まず、学習用データの集合が情報処理装置 1 0 に入力される（Step 1）。次いで、情報処理装置 1 0 は、入力された学習用データの集合を利用し、既に説明した推定機の自動構築方法により推定機を構築する（Step 2）。

【 0 0 6 7 】

次いで、情報処理装置 1 0 は、随時又は所定のタイミングで追加の学習用データを取得する（Step 3）。次いで、情報処理装置 1 0 は、（Step 1）で入力された学習用データの集合に、（Step 3）で取得した学習用データを統合する（Step 4）。このとき、情報処理装置 1 0 は、学習用データのサンプリング処理や重み付け処理を実行し、統合後の学習用データの集合を生成する。そして、情報処理装置 1 0 は、統合後の学習用データの集合を利用し、再び推定機を構築する（Step 2）。このとき、情報処理装置 1 0 は、既に説明した推定機の自動構築方法により推定機を構築する。

【 0 0 6 8 】

また、（Step 2）～（Step 4）の処理は繰り返し実行される。そして、学習用データは、処理が繰り返される度に更新される。例えば、繰り返しの度に学習用データが追加されるようにすれば、推定機の構築処理に利用される学習用データの数が増加するため、推定機の性能が向上する。なお、（Step 4）で実行される学習用データの統合処理においては、情報処理装置 1 0 のリソースをより有効に利用すべく、より有用な学習用データが推定機の構築に利用されるように統合の仕方を工夫する。以下、この工夫について紹介する。

【 0 0 6 9 】

[1 - 3 : 学習用データの統合方法]

学習用データの統合方法について、より詳細に説明する。

【 0 0 7 0 】

（ 1 - 3 - 1 : 特徴量空間における学習用データの分布と推定機の精度 ）

まず、図 1 4 を参照しながら、特徴量空間における学習用データの分布と推定機の精度との関係について考察する。図 1 4 は、特徴量空間における学習用データの分布例を示した説明図である。

【 0 0 7 1 】

1 つの特徴量ベクトルは、1 つの学習用データを構成するデータを基底関数リストに含まれる各基底関数に入力することで得られる。つまり、1 つの学習用データには 1 つの特徴量ベクトル（特徴量座標）が対応する。そのため、特徴量座標の分布を特徴量空間における学習用データの分布と呼ぶことにする。特徴量空間における学習用データの分布は、例えば、図 1 4 のようになる。なお、表現の都合上、図 1 4 の例では 2 次元の特徴量空間

10

20

30

40

50

を考えているが、特徴量空間の次元数はこれに限定されない。

【0072】

さて、図14に例示した特徴量座標の分布を参照すると、第4象限に疎な領域が存在していることに気づくであろう。既に説明した通り、推定関数は、全ての学習用データについて特徴量ベクトルと目的変数との関係がうまく表現されるように回帰/判別学習により生成される。そのため、特徴量座標の密度が疎な領域について、推定関数は、特徴量ベクトルと目的変数との関係をうまく表現できていない可能性が高い。従って、認識処理の対象となる入力データに対応する特徴量座標が上記の疎な領域に位置する場合、高精度の認識結果を期待することは難しい。

【0073】

図15に示すように、学習データセットの数が多くなると疎な領域が生じにくくなり、どの領域に対応する入力データが入力されても高い精度で認識結果を出力することが可能な推定機を構築できるようになると期待される。また、学習データセットの数が比較的少なくても、特徴量座標が特徴量空間において満遍なく分布していれば、高い精度で認識結果を出力することが可能な推定機を構築できるものと期待される。そこで、本件発明者は、学習用データを統合する際に特徴量座標の分布を考慮し、統合後の学習用データの集合に対応する特徴量座標の分布が所定の分布（例えば、一様分布やガウス分布など）となるように調整する仕組みを考案した。

【0074】

(1-3-2: データ統合時にサンプリングする構成)

まず、図16を参照しながら、学習用データをサンプリングする方法について説明する。図16は、学習用データをサンプリングする方法について説明するための説明図である。

【0075】

既に説明したように、オンライン学習を適用する場合、逐次的に学習用データを追加できるため、多量の学習用データを用いて推定機を構築することが可能になる。しかし、情報処理装置10のメモリリソースが限られている場合、学習用データの統合時に、推定機の構築に利用する学習用データの数を絞り込む必要がある。このとき、ランダムに学習用データを間引くのではなく、特徴量座標の分布を考慮して学習用データを間引くことで、推定機の精度を低下させることなく、学習用データの数を絞り込むことができる。例えば、図16に示すように、密な領域に含まれる特徴量座標を多く間引き、疎な領域に含まれる特徴量座標を極力残すようにする。

【0076】

このような方法で学習用データを間引くことにより、統合後の学習用データの集合に対応する特徴量座標の密度が均一になる。つまり、学習用データの数は少なくなったが、特徴量空間の全体に満遍なく特徴量座標が分布しているため、推定関数の生成時に実行する回帰/判別学習の際に特徴量空間の全体が考慮されることになる。その結果、情報処理装置10のメモリリソースが限られていても、高い精度で正しい認識結果を推定することが可能な推定機を構築することが可能になる。

【0077】

(1-3-3: データ統合時に重み付けする構成)

次に、学習用データに重みを設定する方法について説明する。

【0078】

情報処理装置10のメモリリソースが限られている場合、学習用データの統合時に学習用データを間引く方法は有効である。一方、メモリリソースに余裕がある場合、学習用データを間引く代わりに、学習用データに重みを設定することで推定機の性能を向上させることが可能になる。例えば、疎な領域に特徴量座標が含まれる学習用データには大きな重みを設定し、密な領域に特徴量座標が含まれる学習用データには小さな重みを設定する。そして、推定関数の生成時に実行する回帰/判別学習の際に各学習用データに設定された重みを考慮するようにする。

10

20

30

40

50

【 0 0 7 9 】

(1 - 3 - 4 : データ統合時にサンプリング及び重み付けする構成)

また、学習用データをサンプリングする方法と、学習用データに重みを設定する方法とを組み合わせてもよい。例えば、特徴量座標の分布が所定の分布となるように学習用データを間引いた後、間引き後の学習用データの集合に属する学習用データに対し、特徴量座標の密度に応じた重みを設定する。このように、間引き処理と重み付け処理とを組み合わせることにより、メモリリソースが限られていても、より高精度の推定機を構築することが可能になる。

【 0 0 8 0 】

[1 - 4 : 効率的なサンプリング / 重み付け方法]

次に、学習用データの効率的なサンプリング / 重み付け方法について説明する。

【 0 0 8 1 】

(1 - 4 - 1 : サンプリング方法)

まず、図 1 7 を参照しながら、学習用データの効率的なサンプリング方法について説明する。図 1 7 は、学習用データの効率的なサンプリング方法について説明するための説明図である。

【 0 0 8 2 】

図 1 7 に示すように、情報処理装置 1 0 は、全ての学習用データについて特徴量ベクトル (特徴量座標) を算出する (S 2 0 1)。次いで、情報処理装置 1 0 は、算出した特徴量座標を正規化する (S 2 0 2)。例えば、情報処理装置 1 0 は、図 1 8 に示すように、各特徴量について、分散が 1、平均が 0 となるように値を正規化する。

【 0 0 8 3 】

次いで、情報処理装置 1 0 は、ランダムにハッシュ関数 g を生成する (S 2 0 3)。例えば、情報処理装置 1 0 は、下記の式 (1) に示すような 5 ビットの値を出力するハッシュ関数 g を複数生成する。このとき、情報処理装置 1 0 は、 Q 個のハッシュ関数 g_q ($q = 1 \sim Q$) を生成する。但し、関数 h_j ($j = 1 \sim 5$) は、下記の式 (2) により定義される。また、 d 及び $Threshold$ は、乱数により決定される。

【 0 0 8 4 】

但し、特徴量座標の分布を一様分布に近づける場合、 $Threshold$ の決定に用いる乱数として一様乱数を用いる。また、特徴量座標の分布をガウス分布に近づける場合、 $Threshold$ の決定に用いる乱数としてガウス乱数を用いる。他の分布についても同様である。また、 d の決定は、 z_d の算出に用いた基底関数の寄与率に応じた偏りのある乱数を用いて行われる。例えば、 z_d の算出に用いた基底関数の寄与率が大きいほど、 d の発生する確率が高くなる乱数が用いられる。

【 0 0 8 5 】

【数 1】

$$g(Z) = \{h_1(Z), h_2(Z), h_3(Z), h_4(Z), h_5(Z)\}$$

… (1)

$$h_j(Z) = \begin{cases} 1 & (z_d > Threshold) \\ 0 & (z_d \leq Threshold) \end{cases}$$

… (2)

【 0 0 8 6 】

ハッシュ関数 g_q ($q = 1 \sim Q$) を生成すると、情報処理装置 1 0 は、各学習用データに対応する特徴量ベクトル Z をハッシュ関数 g_q に入力し、ハッシュ値を算出する。そし

10

20

30

40

50

て、情報処理装置 10 は、算出したハッシュ値に基づいて学習用データをバケットに割り当てる (S204)。但し、ここで言うバケットとは、ハッシュ値として取り得る値が対応付けられた領域を意味する。

【0087】

例えば、ハッシュ値が 5 ビット、 $Q = 256$ の場合について考えてみよう。この場合、バケットの構成は図 19 のようになる。図 19 に示すように、ハッシュ値が 5 ビットであるから、1 つのハッシュ関数 g_q に対し、32 個のバケット (以下、バケットセット) が設けられる。また、 $Q = 256$ であるから、256 組のバケットセットが設けられる。この例に沿って、学習用データをバケットに割り当てる方法について説明する。

【0088】

ある学習用データに対応する特徴量ベクトル Z が与えられると、256 個のハッシュ関数 $g_1 \sim g_{256}$ を用いて 256 個のハッシュ値が算出される。例えば、 $g_1(Z) = 2$ (10 進数表示) であった場合、情報処理装置 10 は、その学習用データを g_1 に対応するバケットセットの中で 2 に対応するバケットに割り当てる。同様に、 $g_q(Z)$ ($q = 2 \sim 256$) を算出し、各値に対応するバケットに学習用データを割り当てる。図 19 の例では、2 種類の学習用データを白丸と黒丸とで表現し、各バケットとの対応関係を模式的に表現している。

【0089】

このようにして各学習用データをバケットに割り当てると、情報処理装置 10 は、所定の順序でバケットから学習用データを 1 つ選択する (S205)。例えば、情報処理装置 10 は、図 20 に示すように、左上 (ハッシュ関数のインデックス q が小さく、バケットに割り当てられた値が小さい側) から順にバケットを走査し、バケットに割り当てられた学習用データを 1 つ選択する。

【0090】

バケットから学習用データを選択するルールは、図 21 に示した通りである。第 1 に、情報処理装置 10 は、空のバケットをスキップする。第 2 に、情報処理装置 10 は、1 つの学習用データを選択した場合、同じ学習用データを他の全てのバケットから除く。第 3 に、情報処理装置 10 は、1 つのバケットに複数の学習用データが割り当てられている場合にはランダムに 1 つの学習用データを選択する。なお、選択された学習用データの情報は、情報処理装置 10 により保持される。

【0091】

1 つの学習用データを選択した後、情報処理装置 10 は、所定数の学習用データを選択し終えたか否かを判定する (S206)。所定数の学習用データを選択し終えた場合、情報処理装置 10 は、選択した所定数の学習用データを統合後の学習用データの集合として出力し、学習用データの統合に係る一連の処理を終了する。一方、所定数の学習用データを選択し終えていない場合、情報処理装置 10 は、処理をステップ S205 に進める。

【0092】

以上、学習用データの効率的なサンプリング方法について説明した。なお、特徴量空間と上記のバケットとの対応関係は図 22 に示したイメージ図のようになる。また、上記の方法により学習用データのサンプリングを行った結果は、例えば、図 23 (一様分布の例) のようになる。図 23 を参照すると、疎な領域に含まれる特徴量座標は残り、密な領域に含まれる特徴量座標が間引かれていることが分かる。なお、上記のバケットを利用しない場合、学習用データのサンプリングに要する演算負荷は格段に大きくなる点に注意されたい。

【0093】

(1-4-2: 重み付け方法)

次に、図 24 を参照しながら、学習用データの効率的な重み付け方法について説明する。図 24 は、学習用データの効率的な重み付け方法について説明するための説明図である。

【0094】

10

20

30

40

50

図24に示すように、情報処理装置10は、全ての学習用データについて特徴量ベクトル（特徴量座標）を算出する（S211）。次いで、情報処理装置10は、算出した特徴量座標を正規化する（S212）。例えば、情報処理装置10は、図24に示すように、各特徴量について、分散が1、平均が0となるように値を正規化する。

【0095】

次いで、情報処理装置10は、ランダムにハッシュ関数 g を生成する（S213）。例えば、情報処理装置10は、上記の式（1）に示すような5ビットの値を出力するハッシュ関数 g を複数生成する。このとき、情報処理装置10は、 Q 個のハッシュ関数 g_q （ $q = 1 \sim Q$ ）を生成する。但し、関数 h_j （ $j = 1 \sim 5$ ）は、上記の式（2）により定義される。また、 d 及び $Threshold$ は、乱数により決定される。

10

【0096】

但し、特徴量座標の分布を一様分布に近づける場合、 $Threshold$ の決定に用いる乱数として一様乱数を用いる。また、特徴量座標の分布をガウス分布に近づける場合、 $Threshold$ の決定に用いる乱数としてガウス乱数を用いる。他の分布についても同様である。また、 d の決定は、 z_d の算出に用いた基底関数の寄与率に応じた偏りのある乱数を用いて行われる。例えば、 z_d の算出に用いた基底関数の寄与率が大きいほど、 d の発生する確率が高くなる乱数を用いられる。

【0097】

ハッシュ関数 g_q （ $q = 1 \sim Q$ ）を生成すると、情報処理装置10は、各学習用データに対応する特徴量ベクトル Z をハッシュ関数 g_q に入力し、ハッシュ値を算出する。そして、情報処理装置10は、算出したハッシュ値に基づいて学習用データをバケットに割り当てる（S214）。次いで、情報処理装置10は、各学習用データについて密度を算出する（S215）。例えば、図25に示すように、学習データセットがバケットに割り当てられているものとしよう。また、白丸で表現された学習用データに注目する。

20

【0098】

この場合、情報処理装置10は、まず、各ハッシュ関数に対応するバケットセットについて、白丸を含むバケットに割り当てられている学習用データの数をカウントする。例えば、ハッシュ関数 g_1 に対応するバケットセットを参照すると、白丸を含むバケットに割り当てられている学習用データの数は1である。同様に、ハッシュ関数 g_2 に対応するバケットセットを参照すると、白丸を含むバケットに割り当てられている学習用データの数は2である。情報処理装置10は、ハッシュ関数 $g_1 \sim g_{256}$ に対応するバケットセットについて、白丸を含むバケットに割り当てられている学習用データの数をカウントする。

30

【0099】

そして、情報処理装置10は、カウントした数の平均値を算出し、算出した平均値を白丸に対応する学習用データの密度とみなす。同様にして、情報処理装置10は、全ての学習用データの密度を算出する。なお、各学習用データの密度は図26のB図のように表現される。但し、色が濃い部分の密度が高く、色が薄い部分の密度が低い。

【0100】

さて、全ての学習用データについて密度を算出し終わると、情報処理装置10は、処理をステップS217に進める（S216）。ステップS217に進めた場合、情報処理装置10は、算出した密度から各学習用データに設定する重みを算出する（S217）。例えば、情報処理装置10は、密度の逆数を重みに設定する。なお、各学習用データに設定される重みの分布は図27のB図のように表現される。但し、色が濃い部分の重みが大きく、色が薄い部分の重みが小さい。図27を参照すると、密な領域の重みが小さく、疎な領域の重みが大きくなっていることが分かるであろう。

40

【0101】

上記のようにして各学習用データに設定する重みを算出し終わると、情報処理装置10は、重み付けに係る一連の処理を終了する。以上、学習用データの効率的な重み付け方法について説明した。なお、上記のバケットを利用しない場合、学習用データの重み付けに

50

要する演算負荷は格段に大きくなる点に注意されたい。

【 0 1 0 2 】

(1 - 4 - 3 : 組み合わせ方法)

次に、図 2 8 を参照しながら、上記の効率的なサンプリング方法と効率的な重み付け方法とを組み合わせる方法について説明する。図 2 8 は、上記の効率的なサンプリング方法と効率的な重み付け方法とを組み合わせる方法について説明するための説明図である。

【 0 1 0 3 】

図 2 8 に示すように、情報処理装置 1 0 は、まず、学習用データのサンプリング処理を実行する (S 2 2 1)。このサンプリング処理は、図 1 7 に示した処理の流れに沿って実行される。そして、所定数の学習用データが得られると、情報処理装置 1 0 は、得られた学習用データを対象に重み付け処理を実行する (S 2 2 2)。この重み付け処理は、図 2 4 に示した処理の流れに沿って実行される。なお、サンプリング処理の際に算出した特徴量ベクトルやハッシュ関数を流用してもよい。サンプリング処理及び重み付け処理を実行し終わると、情報処理装置 1 0 は、一連の処理を終了する。

10

【 0 1 0 4 】

以上、学習用データの効率的なサンプリング / 重み付け方法について説明した。

【 0 1 0 5 】

[1 - 5 : サンプリング処理及び重み付け処理に関する変形例]

次に、サンプリング処理及び重み付け処理に関する変形例を紹介する。

【 0 1 0 6 】

(1 - 5 - 1 : 変形例 1 (距離に基づく処理))

まず、図 2 9 を参照しながら、特徴量座標間の距離に基づく学習用データのサンプリング方法について説明する。図 2 9 は、特徴量座標間の距離に基づく学習用データのサンプリング方法について説明するための説明図である。

20

【 0 1 0 7 】

図 2 9 に示すように、情報処理装置 1 0 は、まず、ランダムに 1 つの特徴量座標を選択する (S 2 3 1)。次いで、情報処理装置 1 0 は、インデックス j を 1 に初期化する (S 2 3 2)。次いで、情報処理装置 1 0 は、未だ選択されていない J 個の特徴量座標の中から j 番目の特徴量座標を対象座標に設定する (S 2 3 3)。次いで、情報処理装置 1 0 は、既に選択された全ての特徴量座標と対象座標との距離 D を算出する (S 2 3 4)。次いで、情報処理装置 1 0 は、算出した距離 D の最小値 D_{min} を抽出する (S 2 3 5)。

30

【 0 1 0 8 】

次いで、情報処理装置 1 0 は、 $j = J$ であるか否かを判定する (S 2 3 6)。 $j = J$ である場合、情報処理装置 1 0 は、処理をステップ S 2 3 7 に進める。一方、 $j < J$ である場合、情報処理装置 1 0 は、処理をステップ S 2 3 3 に進める。処理をステップ S 2 3 7 に進めた場合、情報処理装置 1 0 は、最小値 D_{min} が最大となる対象座標 (特徴量座標) を選択する (S 2 3 7)。次いで、情報処理装置 1 0 は、ステップ S 2 3 1 及び S 2 3 7 において選択された特徴量座標の数が所定数に達したか否かを判定する (S 2 3 8)。

【 0 1 0 9 】

ステップ S 2 3 1 及び S 2 3 7 において選択された特徴量座標の数が所定数に達した場合、情報処理装置 1 0 は、選択された特徴量座標に対応する学習用データを統合後の学習用データの集合として出力し、一連の処理を終了する。一方、ステップ S 2 3 1 及び S 2 3 7 において選択された特徴量座標の数が所定数に達していない場合、情報処理装置 1 0 は、処理をステップ S 2 3 2 に進める。

40

【 0 1 1 0 】

以上、特徴量座標間の距離に基づく学習用データのサンプリング方法について説明した。

【 0 1 1 1 】

(1 - 5 - 2 : 変形例 2 (クラスタリングに基づく処理))

次に、クラスタリングに基づく学習用データのサンプリング / 重み付け方法について説

50

明する。なお、以下ではサンプリング方法及び重み付け方法についてそれぞれ別々に説明するが、これらの方法を組み合わせてもよい。

【0112】

(データセットの選択)

まず、図30を参照しながら、クラスタリングに基づく学習用データのサンプリング方法について説明する。図30は、クラスタリングに基づく学習用データのサンプリング方法について説明するための説明図である。

【0113】

図30に示すように、まず、情報処理装置10は、特徴量ベクトルを所定数のクラスタに分類する(S241)。クラスタリング手法としては、例えば、k-means法や階層的クラスタリングなどの手法が利用可能である。次いで、情報処理装置10は、各クラスタから順に1つずつ特徴量ベクトルを選択する(S242)。そして、情報処理装置10は、選択した特徴量ベクトルに対応する学習用データを統合後の学習用データの集合として出力し、一連の処理を終了する。

10

【0114】

(重みの設定)

次に、図31を参照しながら、クラスタリングに基づく学習用データの重み付け方法について説明する。図31は、クラスタリングに基づく学習用データの重み付け方法について説明するための説明図である。

【0115】

図31に示すように、まず、情報処理装置10は、特徴量ベクトルを所定数のクラスタに分類する(S251)。クラスタリング手法としては、例えば、k-means法や階層的クラスタリングなどの手法が利用可能である。次いで、情報処理装置10は、各クラスタの要素数をカウントし、要素数の逆数を算出する(S252)。そして、情報処理装置10は、算出した要素数の逆数を重みとして出力し、一連の処理を終了する。

20

【0116】

以上、クラスタリングに基づく学習用データのサンプリング/重み付け方法について説明した。

【0117】

(1-5-3:変形例3(密度推定手法に基づく処理))

次に、密度推定手法に基づく学習用データのサンプリング/重み付け方法について説明する。なお、以下ではサンプリング方法及び重み付け方法についてそれぞれ別々に説明するが、これらの方法を組み合わせてもよい。

30

【0118】

(データセットの選択)

まず、図32を参照しながら、密度推定手法に基づく学習用データのサンプリング方法について説明する。図32は、密度推定手法に基づく学習用データのサンプリング方法について説明するための説明図である。

【0119】

図32に示すように、まず、情報処理装置10は、特徴量座標の密度をモデル化する(S261)。密度のモデル化には、例えば、GMM(Gaussian Mixture Model)などの密度推定手法が利用される。次いで、情報処理装置10は、構築したモデルに基づいて各特徴量座標の密度を算出する(S262)。次いで、情報処理装置10は、未だ選択されていない特徴量座標の中から、密度の逆数に比例する確率でランダムに特徴量座標を選択する(S263)。

40

【0120】

次いで、情報処理装置10は、所定数の特徴量座標を選択したか否かを判定する(S264)。所定数の特徴量座標を選択していない場合、情報処理装置10は、処理をステップS263に進める。一方、所定数の特徴量座標を選択した場合、情報処理装置10は、選択した特徴量座標に対応する学習用データを統合後の学習用データの集合として出力し

50

、一連の処理を終了する。

【0121】

(重みの設定)

次に、図33を参照しながら、密度推定手法に基づく学習用データの重み付け方法について説明する。図33は、密度推定手法に基づく学習用データの重み付け方法について説明するための説明図である。

【0122】

図33に示すように、まず、情報処理装置10は、特徴量座標の密度をモデル化する(S271)。密度のモデル化には、例えば、GMMなどの密度推定手法が利用される。次いで、情報処理装置10は、構築したモデルに基づいて各特徴量座標の密度を算出する(S272)。そして、情報処理装置10は、算出した密度の逆数を重みに設定し、一連の処理を終了する。

10

【0123】

以上、密度推定手法に基づく学習用データのサンプリング/重み付け方法について説明した。

【0124】

以上、後述する実施形態において利用可能な基盤技術について説明した。但し、後述する実施形態に係る技術は、ここで説明した基盤技術の全てを利用しなくてもよいし、当該基盤技術を変形して利用したり、或いは、他の機械学習アルゴリズムを組み合わせ利用したりしてもよい点に注意されたい。

20

【0125】

<2:実施形態>

以下、本技術の一実施形態について説明する。

【0126】

[2-1:思考ルーチンの自動構築方法]

本実施形態は、ロボットのようなエージェントの思考ルーチンや様々なゲームに登場するNPC(Non-Player Character)の思考ルーチンを自動構築する技術に関する。例えば、本実施形態は、ある状態SにおかれたNPCが次にとる行動aを決定する思考ルーチンを自動構築する技術に関する。本稿においては、状態Sの入力に応じて行動aを出力するプログラムを思考ルーチンと呼ぶことにする。また、以下では、NPCの行動aを決定する思考ルーチンを例に挙げて説明を進めることにする。もちろん、ロボットなどの行動を決定する思考ルーチンも同様に自動構築することが可能である。

30

【0127】

(2-1-1:思考ルーチンとは)

上記の通り、本稿に言う思考ルーチンは、図34に示すように、状態Sの入力に応じて行動aを出力するプログラムである。なお、状態Sとは、ある瞬間に、行動aを決定すべきNPCがおかれた環境を意味する。例えば、図34に示すように、2つのNPC(NPC#1及び#2)が対戦する対戦ゲームについて考えてみよう。この対戦ゲームは、NPC#1及び#2がそれぞれヒットポイントを有しており、ダメージを受けるとヒットポイントが減少していく仕組みになっているものとする。この例において、ある瞬間における状態Sは、NPC#1及び#2のヒットポイント及び位置関係になる。

40

【0128】

この状態Sが入力されると、思考ルーチンは、NPC#1がNPC#2にダメージを与え、最終的にNPC#2のヒットポイントを0にできることが期待されるNPC#1の行動aを決定する。例えば、NPC#1のヒットポイントが十分に残っており、NPC#2のヒットポイントが僅かである場合、思考ルーチンは、NPC#1が多少のダメージを受けることを許容してNPC#2に素早くダメージを与える行動aを決定するかもしれない。また、NPC#1のヒットポイントが残り僅かであり、NPC#2のヒットポイントが十分に残っている場合、思考ルーチンは、NPC#1がダメージを受けないようにしつつ、NPC#2にダメージを与える行動aを決定するだろう。

50

【 0 1 2 9 】

これまで、NPCの行動を決定する思考ルーチンは、熟練した技術者により長い時間をかけて設計されていた。もちろん、NPCの行動をランダムに決定する思考ルーチンも存在するであろう。しかし、賢いNPCの行動を実現することが可能な思考ルーチンを構築するには、ユーザ操作の分析や環境に応じた最適な行動の研究が欠かせなかった。さらに、こうした分析や研究の結果を踏まえて、環境に応じたNPCの最適な行動を決定するための条件設計を行う必要があった。そのため、思考ルーチンの構築には長い時間と大きな労力とが必要であった。こうした事情を踏まえ、本件発明者は、このような思考ルーチンを人手に依らずに自動構築する技術を開発した。

【 0 1 3 0 】

(2 - 1 - 2 : 思考ルーチンの構成)

図35に示すように、本実施形態に係る思考ルーチンは、行動履歴データに基づく思考ルーチンの自動構築技術により生成される。この行動履歴データは、状態 S 、行動 a 、報酬 r により構成される。例えば、状態 $S = S_1$ において、NPC # 1 が行動 $a =$ “ 右へ移動 ” をとった場合にNPC # 2 からダメージを受けてヒットポイントが0になったとしよう。この場合、行動履歴データは、状態 $S = S_1$ 、行動 $a =$ “ 右へ移動 ”、報酬 $r =$ “ 0 ” となる。このような構成を有する行動履歴データを予め蓄積しておき、この行動履歴データを学習データとする機械学習により思考ルーチンを自動構築することができる。

【 0 1 3 1 】

本実施形態に係る思考ルーチンは、図36に示すような構成を有する。図36に示すように、この思考ルーチンは、状態 S の入力に応じてNPCがとりうる行動 a をリストアップし、各行動 a についてNPCが得るであろう報酬 r の推定値(以下、推定報酬 y)を算出する。そして、思考ルーチンは、推定報酬 y が最も高い行動 a を選択する。なお、推定報酬 y は、報酬推定機を利用して算出される。この推定報酬機は、状態 S 及び行動 a の入力に応じて推定報酬 y を出力するアルゴリズムである。また、この報酬推定機は、行動履歴データを学習データとする機械学習により自動構築される。例えば、先に紹介した推定機の自動構築方法を応用することにより、報酬推定機を自動構築することができる。

【 0 1 3 2 】

報酬推定機は、図37に示すように、基底関数リスト(ϕ_1, \dots, ϕ_M)及び推定関数 f により構成される。基底関数リスト(ϕ_1, \dots, ϕ_M)は、 M 個の基底関数 ϕ_k ($k = 1 \sim M$)を含む。また、基底関数 ϕ_k は、入力データ X (状態 S 及び行動 a)の入力に応じて特徴量 z_k を出力する関数である。さらに、推定関数 f は、 M 個の特徴量 z_k ($k = 1 \sim M$)を要素として含む特徴量ベクトル $Z = (z_1, \dots, z_M)$ の入力に応じて推定報酬 y を出力する関数である。基底関数 ϕ_k は、予め用意された1又は複数の処理関数を組み合わせることで生成される。処理関数としては、例えば、三角関数、指数関数、四則演算、デジタルフィルタ、微分演算、中央値フィルタ、正規化演算などが利用可能である。

【 0 1 3 3 】

また、本実施形態に係る思考ルーチンの自動構築技術は、自動構築された思考ルーチンを利用してNPCを行動させ、その行動の結果として得られた行動履歴データを追加した行動履歴データを利用して思考ルーチンを更新する。但し、行動履歴データの追加には、先に紹介したオンライン学習に係る技術を利用することができる。

【 0 1 3 4 】

(2 - 1 - 3 : 報酬推定機の構築方法)

例えば、オンライン学習に係る技術を利用した報酬推定機の構築及び更新は、図38に示すような処理の流れに沿って行われる。なお、これらの処理は、情報処理装置10により実行されるものとする。図38に示すように、まず、行動履歴データが情報処理装置10に入力される(Step 1)。

【 0 1 3 5 】

(Step 1)において、情報処理装置10は、予め設計された簡易な思考ルーチンを用いて行動 a を決定しながらNPCを環境中で振る舞わせ、行動履歴データ($S, a,$

10

20

30

40

50

r)を得る。この簡易な思考ルーチンは、強化学習の分野においてInnate(赤ちゃんが行う本能的な動きに相当)と呼ばれる。このInnateは、NPCが取り得るアクションの中からランダムに行動を選択するものであってもよい。この場合、Innateの設計も不要になる。情報処理装置10は、所定数の行動履歴データが得られるまでInnateに基づくNPCの行動を繰り返し実行する。次いで、情報処理装置10は、入力された行動履歴データを利用し、既に説明した推定機の自動構築方法と同様にして報酬推定機を構築する(Step 2)。

【0136】

次いで、情報処理装置10は、随時又は所定のタイミングで追加の行動履歴データを取得する(Step 3)。次いで、情報処理装置10は、(Step 1)で入力された行動履歴データと、(Step 3)で取得した行動履歴データとを統合する(Step 4)。このとき、情報処理装置10は、行動履歴データのサンプリング処理や重み付け処理を実行し、統合後の行動履歴データを生成する。そして、情報処理装置10は、統合後の行動履歴データを利用し、再び報酬推定機を構築する(Step 2)。また、(Step 2)~(Step 4)の処理は繰り返し実行される。そして、行動履歴データは、処理が繰り返される度に更新される。

【0137】

以上、思考ルーチンの自動構築方法について簡単に説明した。ここではNPCの行動を決定する思考ルーチンの自動構築方法について述べたが、行動履歴データの構成を変えることで様々な種類の思考ルーチンを同じように自動構築することができる。つまり、本実施形態の技術を適用することにより、統一的な仕組みで様々な思考ルーチンを構築できるようになる。また、自動的に思考ルーチンが構築されるため、思考ルーチンの構築に人が時間を費やさずに済み、労力が大幅に軽減される。

【0138】

[2-2: 情報処理装置10の構成]

ここで、図39及び図40を参照しながら、本実施形態に係る情報処理装置10の機能構成について説明する。図39は、本実施形態に係る情報処理装置10の全体的な機能構成を示した説明図である。一方、図40は、本実施形態に係る情報処理装置10を構成する報酬推定機構築部12の詳細な機能構成を示した説明図である。

【0139】

(全体的な機能構成)

まず、図39を参照しながら、全体的な機能構成について説明する。図39に示すように、情報処理装置10は、主に、行動履歴データ取得部11と、報酬推定機構築部12と、入力データ取得部13と、行動選択部14とにより構成される。

【0140】

思考ルーチンの構築処理が開始されると、行動履歴データ取得部11は、報酬推定機の構築に利用する行動履歴データを取得する。例えば、行動履歴データ取得部11は、簡易な思考ルーチン(Innate)に基づいて繰り返しNPCを行動させ、所定数の行動履歴データを取得する。但し、行動履歴データ取得部11は、記憶装置(非図示)に予め格納された行動履歴データを読み出したり、或いは、行動履歴データを提供するシステムなどからネットワークを介して行動履歴データを取得したりしてもよい。

【0141】

行動履歴データ取得部11により取得された行動履歴データは、報酬推定機構築部12に入力される。行動履歴データが入力されると、報酬推定機構築部12は、入力された行動履歴データに基づく機械学習により報酬推定機を構築する。例えば、報酬推定機構築部12は、既に説明した遺伝アルゴリズムに基づく推定機の自動構築方法を利用して報酬推定機を構築する。また、行動履歴データ取得部11から追加の行動履歴データが入力された場合、報酬推定機構築部12は、行動履歴データを統合し、統合後の行動履歴データを利用して報酬推定機を構築する。

【0142】

10

20

30

40

50

報酬推定機構築部 1 2 により構築された報酬推定機は、行動選択部 1 4 に入力される。この報酬推定機は、任意の入力データ（状態 S）に対して最適な行動を選択するために利用される。入力データ取得部 1 3 により入力データ（状態 S）が取得されると、取得された入力データは、行動選択部 1 4 に入力される。入力データが入力されると、行動選択部 1 4 は、入力された入力データが示す状態 S 及び状態 S において NPC がとりうる行動 a を報酬推定機に入力し、報酬推定機から出力される推定報酬 y に基づいて行動 a を選択する。例えば、図 3 6 に示すように、行動選択部 1 4 は、推定報酬 y が最も高くなる行動 a を選択する。

【 0 1 4 3 】

以上、情報処理装置 1 0 の全体的な機能構成について説明した。

10

【 0 1 4 4 】

（報酬推定機構築部 1 2 の機能構成）

次に、図 4 0 を参照しながら、報酬推定機構築部 1 2 の機能構成について詳細に説明する。図 4 0 に示すように、報酬推定機構築部 1 2 は、基底関数リスト生成部 1 2 1 と、特徴量計算部 1 2 2 と、推定関数生成部 1 2 3 と、行動履歴データ統合部 1 2 4 とにより構成される。

【 0 1 4 5 】

思考ルーチンの構築処理が開始されると、まず、基底関数リスト生成部 1 2 1 は、基底関数リストを生成する。そして、基底関数リスト生成部 1 2 1 により生成された基底関数リストは、特徴量計算部 1 2 2 に入力される。また、特徴量計算部 1 2 2 には、行動履歴データが入力される。基底関数リスト及び行動履歴データが入力されると、特徴量計算部 1 2 2 は、入力された行動履歴データを基底関数リストに含まれる各基底関数に入力して特徴量を算出する。特徴量計算部 1 2 2 により算出された特徴量の組（特徴量ベクトル）は、推定関数生成部 1 2 3 に入力される。

20

【 0 1 4 6 】

特徴量ベクトルが入力されると、推定関数生成部 1 2 3 は、入力された特徴量ベクトル及び行動履歴データを構成する報酬値 r に基づいて回帰 / 判別学習により推定関数を生成する。なお、遺伝アルゴリズムに基づく推定機の構築方法を適用する場合、推定関数生成部 1 2 3 は、生成した推定関数に対する各基底関数の寄与率（評価値）を算出し、その寄与率に基づいて終了条件を満たすか否かを判定する。終了条件を満たす場合、推定関数生成部 1 2 3 は、基底関数リスト及び推定関数を含む報酬推定機を出力する。

30

【 0 1 4 7 】

一方、終了条件を満たさない場合、推定関数生成部 1 2 3 は、生成した推定関数に対する各基底関数の寄与率を基底関数リスト生成部 1 2 1 に通知する。この通知を受けた基底関数リスト生成部 1 2 1 は、遺伝アルゴリズムにより各基底関数の寄与率に基づいて基底関数リストを更新する。基底関数リストを更新した場合、基底関数リスト生成部 1 2 1 は、更新後の基底関数リストを特徴量計算部 1 2 2 に入力する。更新後の基底関数リストが入力された場合、特徴量計算部 1 2 2 は、更新後の基底関数リストを用いて特徴量ベクトルを算出する。そして、特徴量計算部 1 2 2 により算出された特徴量ベクトルは、推定関数生成部 1 2 3 に入力される。

40

【 0 1 4 8 】

上記のように、遺伝アルゴリズムに基づく推定機の構築方法を適用する場合、終了条件が満たされるまで、推定関数生成部 1 2 3 による推定関数の生成処理、基底関数リスト生成部 1 2 1 による基底関数リストの更新処理、及び特徴量計算部 1 2 2 による特徴量ベクトルの算出処理が繰り返し実行される。そして、終了条件が満たされた場合、推定関数生成部 1 2 3 から報酬推定機が出力される。

【 0 1 4 9 】

また、追加の行動履歴データが入力されると、入力された追加の行動履歴データは、特徴量計算部 1 2 2 及び行動履歴データ統合部 1 2 4 に入力される。追加の行動履歴データが入力されると、特徴量計算部 1 2 2 は、追加の行動履歴データを基底関数リストに含

50

れる各基底関数に入力して特徴量を生成する。そして、追加の行動履歴データに対応する特徴量ベクトル及び既存の行動履歴データに対応する特徴量ベクトルは、行動履歴データ統合部 124 に入力される。なお、行動履歴データ統合部 124 には、既存の行動履歴データも入力されているものとする。

【0150】

行動履歴データ統合部 124 は、先に紹介したデータセットの統合方法を応用して既存の行動履歴データと追加の行動履歴データとを統合する。例えば、行動履歴データ統合部 124 は、特徴量空間において特徴量ベクトルにより示される座標（特徴量座標）の分布が所定の分布となるように行動履歴データを間引いたり、行動履歴データに重みを設定したりする。行動履歴データを間引いた場合、間引き後の行動履歴データが統合後の行動履歴データとして利用される。一方、行動履歴データに重みを設定した場合、推定関数生成部 123 による回帰 / 判別学習の際に各行動履歴データに設定された重みが考慮される。

10

【0151】

行動履歴データが統合されると、統合後の行動履歴データを用いて報酬推定機の自動構築処理が実行される。具体的には、行動履歴データ統合部 124 から推定関数生成部 123 に統合後の行動履歴データと、統合後の行動履歴データに対応する特徴量ベクトルとが入力され、推定関数生成部 123 により推定関数が生成される。また、遺伝アルゴリズムに基づく推定機の構築方法を適用する場合、統合後の行動履歴データを利用して推定関数の生成、寄与率の算出、基底関数リストの更新などの処理が実行される。

【0152】

以上、報酬推定機構築部 12 の詳細な機能構成について説明した。

20

【0153】

以上、本実施形態に係る情報処理装置 10 の構成について説明した。上記の構成を適用することにより、任意の状態 S から NPC がとるべき次の行動 a を決定する思考ルーチンを自動構築することができる。また、この思考ルーチンを利用して賢く NPC を行動させることが可能になる。なお、利用する行動履歴データを変えることで、ロボットなどのエージェントについても同様に賢く行動させることが可能になる。

【0154】

[2-3：効率的な推定報酬機の構築方法]

これまで、先に紹介した推定機の自動構築方法に基づく思考ルーチンの自動構築方法について説明してきた。確かに、この方法を適用すると、思考ルーチンを自動構築することが可能になる。しかし、賢く行動する NPC の思考ルーチンを自動構築するには、ある程度長い時間をかけて学習処理を繰り返し実行する必要がある。そこで、本件発明者は、より効率良く高性能な推定報酬機を構築する方法を考案した。

30

【0155】

以下、図 4 1 及び図 4 2 を参照しながら、効率的な推定報酬機の構築方法について説明する。この方法は、より学習効率の高い行動履歴データを取得する方法に関する。より学習効率の高い行動履歴データとは、より推定報酬が高く、より推定誤差が大きく、かつ、特徴量空間における密度が疎な領域にある特徴量座標に対応するデータである。そこで、図 4 2 に示す 3 つのスコアを導入する。1 つ目は、推定報酬が高いほど大きな値となる報酬スコアである。2 つ目は、特徴量空間における密度が疎であるほど大きな値となる未知スコアである。3 つ目は、推定誤差が大きいほど大きな値となる誤差スコアである。

40

【0156】

例えば、図 4 1 に示した行動 a_1 、 a_2 、 a_3 に注目しよう。仮に、鎖線で囲まれた領域は、推定誤差の小さい領域であるとする。また、図の右上方向に向かうにつれて推定報酬が高くなっていると。この場合、行動 a_1 は、報酬スコアが比較的高く、未知スコアが比較的高く、誤差スコアが比較的低い行動であると言える。また、行動 a_2 は、報酬スコアが比較的低く、未知スコアが比較的低く、誤差スコアが比較的高い行動であると言える。そして、行動 a_3 は、報酬スコアが比較的高く、未知スコアが比較的高く、誤差スコアが比較的高い行動であると言える。

50

【0157】

より報酬スコアの高い行動を優先的に選択することにより、高い報酬を実現するために必要な行動履歴データを収集することができる。また、より未知スコアが高いか、より誤差スコアが高い行動を優先的に選択することにより、その行動を選択した結果が不定であるような行動履歴データを収集することができる。例えば、図41の例では、行動 a_3 を選択することにより、より高い報酬を得られることが期待され、かつ、その行動を選択した結果が不定であるような行動履歴データを収集できると考えられる。図38に示した処理のうち、(Step 1)及び/又は(Step 3)において上記の方法による行動履歴データの取得を行うことで、(Step 2)における報酬推定機の構築をより効率的に実現することが可能になる。

10

【0158】

以上、効率的な推定報酬機の構築方法について説明した。

【0159】

[2-4: (変形例1)アクションスコア推定機を用いる思考ルーチン]

さて、これまでは報酬推定機を用いて報酬を推定し、推定した報酬に基づいて行動を選択する思考ルーチンについて考えてきた。ここでは、図44に示すように、アクションスコア推定機を用いてアクションスコアを推定し、推定したアクションスコアに基づいて行動を選択する思考ルーチンについて考えてみたい。ここで言うアクションスコアとは、とりうる各行動に対応付けられたスコアであり、対応する行動をとることで好ましい結果が得られる確率の高さを表す。

20

【0160】

アクションスコアを利用する場合、行動履歴データは、図43に示すような形で与えられる。まず、情報処理装置10は、これまで説明してきた行動履歴データと同様にして状態 S 、行動 a 、報酬 r の組を収集する。その後、情報処理装置10は、報酬 r に基づいてアクションスコアを算出する。

【0161】

例えば、状態 $S = S_1$ において行動 $a = "R (右へ移動)"$ をとった場合に報酬 $r = "0"$ が得られたものとしよう。この場合、行動 $a = "R"$ に対応するアクションスコアは $"0"$ となり、それ以外の行動 $("L" "N" "J")$ に対応するアクションスコアは $"1"$ となる。その結果、状態 $S = S_1$ 及び行動 $a = "R"$ に対応するアクションスコア $(R, L, N, J) = (0, 1, 1, 1)$ が得られる。

30

【0162】

また、状態 $S = S_2$ において行動 $a = "L (左へ移動)"$ をとった場合に報酬 $r = "1"$ が得られたものとしよう。この場合、行動 $a = "L"$ に対応するアクションスコアは $"1"$ となり、それ以外の行動 $("R" "N" "J")$ に対応するアクションスコアは $"0"$ となる。その結果、状態 $S = S_2$ 及び行動 $a = "L"$ に対応するアクションスコア $(R, L, N, J) = (0, 1, 0, 0)$ が得られる。

【0163】

上記のようにして得られた状態 S 、行動 a 、アクションスコアの組を行動履歴データとして利用すると、機械学習により、状態 S の入力に応じてアクションスコアの推定値を算出するアクションスコア推定機が得られる。例えば、遺伝アルゴリズムに基づく推定機の自動構築方法を適用すれば、高性能なアクションスコア推定機を自動構築することができる。また、行動履歴データを収集する際に、効率的な報酬推定機の構築方法と同様の方法を用いれば、効率的にアクションスコア推定機を自動構築することができる。

40

【0164】

アクションスコア推定機を用いる場合、思考ルーチンの構成は図44のようになる。つまり、状態 S を思考ルーチンに入力すると、思考ルーチンは、アクションスコア推定機に状態 S を入力し、アクションスコアの推定値を算出する。そして、思考ルーチンは、アクションスコアの推定値が最も高い行動を選択する。例えば、図44に示すように、アクションスコアの推定値が $(R, L, J, N) = (0.6, 0.3, 0.4, 0.2)$ であっ

50

た場合、思考ルーチンは、推定値“0.6”に対応する行動“R”を選択する。

【0165】

以上、アクションスコア推定機を用いる思考ルーチンについて説明した。

【0166】

[2-5：(変形例2)予測機を用いた報酬の推定]

次に、予測機を用いた報酬の推定方法について説明する。なお、ここで言う予測機とは、ある時刻 t_1 における状態 $S(t_1)$ 及び状態 $S(t_1)$ においてNPCがとった行動 $a(t_1)$ を入力した場合に、次の時刻 t_2 における状態 $S(t_2)$ を出力するアルゴリズムのことを意味する。

【0167】

(2-5-1：予測機の構築方法)

上記の予測機は、図45に示すような方法で構築される。図45に示すように、時刻毎に取得された行動履歴データが学習データとして利用される。例えば、時刻 t_2 において状態 S_2 にあるNPCが何もしなかった場合に好ましい結果が得られた場合、行動履歴データは、時刻 $t = t_2$ 、状態 $S = S_2$ 、行動 $a =$ “何もせず”、報酬 $r =$ “1”となる。なお、予測機の自動構築方法については、特願2009-277084号明細書に詳しく記載されている。同明細書には、ある時点までの観測値から将来の時点における観測値を予測する予測機を機械学習により自動構築する方法が記載されている。

【0168】

(2-5-2：報酬の推定方法)

上記の予測機を利用すると、図46に示すように、将来得るであろう報酬を推定することが可能になる。例えば、時刻 t において状態 $S(t)$ にあるNPCが行動 $a(t)$ をとった場合に時刻 $t+1$ において実現される状態 $S(t+1)$ を予測し、その状態 $S(t+1)$ においてNPCがとりうる行動毎に推定報酬 y を算出することができるようになる。そのため、時刻 $t+1$ において推定される報酬に基づいて時刻 t においてNPCがとるべき行動を選択することができるようになる。また、図47に示すように、予測機を繰り返し用いて数ステップ先の状態 $S(t+q)$ から推定される推定報酬 y を算出することもできる。この場合、各時刻においてNPCがとりうる行動の組み合わせを考慮し、最終的に最も高い推定報酬が得られる行動の組み合わせを選択することができるようになる。

【0169】

以上、予測機を用いた報酬の推定方法について説明した。

【0170】

[2-6：(変形例3)複数エージェントの同時学習]

さて、これまでは1つのNPCに注目して最適な行動を選択する思考ルーチンの構築方法について考えてきた。しかし、2つ以上のNPCがとる行動を同時に考慮して思考ルーチンを構築することも可能である。2つのNPCが同じ環境中で行動する場合、両NPCがとる行動は状態 S に反映される。そのため、この方法を適用すると、他のNPCが最も高い推定報酬を見込める行動を選択して行動する環境中において、自身のNPCが最も高い推定報酬を見込める行動を選択するような思考ルーチンを自動構築することができる。例えば、MinMax法などを用いることにより、このような思考ルーチンの自動構築が実現される。以上、複数エージェントの同時学習について説明した。

【0171】

以上、本技術の一実施形態について説明した。

【0172】

<3：応用例>

次に、本実施形態の技術を具体的に応用する方法について紹介する。

【0173】

[3-1：「三目並べ」への応用]

まず、図48～図54を参照しながら、本実施形態に係る技術を「三目並べ」へ応用する方法について説明する。図48に示すように、「三目並べ」の主なルールは、(1)交

10

20

30

40

50

互に手を打つ、(2)先に3つのマークが1列に並んだ方が勝ち、の2点である。また、「三目並べ」において、状態Sは盤面であり、行動aは各プレイヤーが打つ手である。

【0174】

「三目並べ」は、互いに最適な手を打つと必ず引き分けになることが知られている。このような完全情報ゲームに用いられる思考ルーチンの多くは、静的評価関数と先読みアルゴリズムとにより構成されている。この静的評価関数は、ある局面の有利/不利を数値化する関数である。例えば、図49に示すような局面が与えられた場合、静的評価関数は、その局面の有利/不利を表す数値y(“不利”:-1、“どちらでもない”:0、“有利”:+1など)を出力する。本実施形態の場合、この静的評価関数の機能は、報酬推定機により実現される。

10

【0175】

また、先読みアルゴリズムは、先の手を読み、将来の静的評価関数の出力値がより高くなるような手を選択するアルゴリズムである。例えば、先読みアルゴリズムは、Minimax法などを利用して実現される。例えば、図50に示すように、先読みアルゴリズムは、自分の手番で手を打った後に、相手の手番で相手が打つ可能性のある手を想定し、想定した各手に対して自分が打てる手を想定して、自分が最も有利になる手を選択する。

【0176】

ところで、上記のような静的評価関数は、これまで人手により設計されていた。例えば、将棋のAIとして有名なボナンザでさえ、静的評価関数で考慮する局面の特徴などの設計事項は人手により設計されていた。また、ゲームの種類が変わると、特徴量の設計も変更する必要がある。そのため、これまでは試行錯誤を繰り返しながら静的評価関数をゲーム毎に人手で設計する必要があった。しかし、本実施形態に係る技術を適用すると、人手による設計作業を省いて思考ルーチンを自動構築することが可能になる。

20

【0177】

「三目並べ」の場合、図51に示すように、状態S及び行動aを3×3のマトリックスで表現する。但し、状態Sは、自分の手番となつた時点の盤面を表す。また、自分の手番で打った手を反映した盤面を(S, a)と表現する。さらに、自分の手を“1”、相手の手を“-1”、空白を“0”と表現する。つまり、盤面及び手を数値で表現する。このようにして盤面及び手が数値で表現できると、本実施形態に係る報酬推定機の自動構築方法を用いて思考ルーチンを自動構築することが可能になる。

30

【0178】

例えば、情報処理装置10は、まず、ランダムな場所に自分の手と相手の手とを打つInnateを利用して行動履歴データを生成する。上記の通り、(S, a)は、3×3マトリックスにより表現される。また、情報処理装置10は、図52に示すように、勝ちに至るまでに打った全ての手に対応する(S, a)に報酬“1”を与える。一方、情報処理装置10は、図53に示すように、負けに至るまでに打った全ての手に対応する(S, a)に報酬“-1”を与える。このようにして行動履歴データを蓄積すると、情報処理装置10は、蓄積した行動履歴データを利用して報酬推定機を構築する。

【0179】

実際に手を選択する場合、情報処理装置10は、図54に示すように、報酬推定機を利用して現在の状態Sから推定報酬yを算出し、推定報酬yが最大となる手を選択する。図54の例では、最大の推定報酬に対応する手(C)が選択される。なお、図54の例では1手先の報酬を評価して手の選択を行っているが、対戦相手についても同じように推定報酬を算出し、Minimax法などを用いて数手先読みした結果を用いて現在の手を選択するように構成してもよい。

40

【0180】

また、学習により得られた報酬推定機を用いて常に最適な行動を選択するように構成すると、NPCによる手の選択が毎回同じになってしまうことがある。そこで、推定報酬を算出する工程に何らかのランダムネスを加えてもよい。例えば、報酬推定機により算出した推定報酬に僅かだけ乱数を加える方法が考えられる。また、遺伝アルゴリズムに基づく

50

機械学習により報酬推定機を算出している場合、学習世代毎に算出される報酬推定機を保持しておき、利用する報酬推定機をランダムに切り替えるように構成してもよい。

【0181】

以上、「三目並べ」への応用について説明した。

【0182】

[3-2:「対戦ゲーム」への応用]

次に、図55～図62を参照しながら、本実施形態に係る技術を「対戦ゲーム」へと応用する方法について説明する。ここで考える「対戦ゲーム」の主なルールは、図55に示すように、(1)2人対戦ゲームであること、(2)各プレイヤーの行動は「左移動」「右移動」「左右移動なし」「ジャンプ」「ジャンプなし」の組み合わせであること、(3)相手のプレイヤーを踏んだらY軸方向の加速度差に応じて相手にダメージを与えられること、の3点である。また、ヒットポイントが0になったプレイヤーが負けである。なお、「対戦ゲーム」への応用には、先に説明したアクションスコア推定機を用いる思考ルーチンの構築方法が用いられる。

【0183】

この場合、状態Sとしては、自分の絶対座標、相手の絶対座標、時刻が利用される。そのため、状態Sは、図56に示すように3次元マトリックスにより表現される。また、ここでは、3次元マトリックスで表現される状態Sの入力に応じて5つの要素(N, L, R, J, NJ)を持つアクションスコアを推定するアクションスコア推定機の自動構築方法について考える。但し、要素Nは、行動a = “左右移動なし”に対応するアクションスコアである。また、要素Lは、行動a = “左移動”に対応するアクションスコアである。要素Rは、行動a = “右移動”に対応するアクションスコアである。要素Jは、行動a = “ジャンプ”に対応するアクションスコアである。要素NJは、行動a = “ジャンプなし”に対応するアクションスコアである。

【0184】

行動履歴データを収集するためのInnateとしては、例えば、完全にランダムにプレイヤーの行動を選択するものが用いられる。例えば、このInnateは、N(左右移動なし)、L(左移動)、R(右移動)の中から1つの行動をランダムに選び、選んだ行動に組み合わせる行動をJ(ジャンプ)又はNJ(ジャンプなし)からランダムに1つ選ぶ。また、情報処理装置10は、図57に示すように、自分が相手にダメージを与えた時点で、前回自分又は相手がダメージを受けた時点から現時点までの行動履歴データの報酬を1に設定する。一方、自分が相手からダメージを受けた場合、情報処理装置10は、図57に示すように、前回自分又は相手がダメージを受けた時点から現時点までの行動履歴データの報酬を0に設定する。

【0185】

なお、報酬が1に設定された行動履歴データについて、情報処理装置10は、実際に行った行動のアクションスコアを1、行わなかった行動のアクションスコアを0に設定する。一方、報酬が0に設定された行動履歴データについて、情報処理装置10は、実際に行った行動のアクションスコアを0に設定し、行わなかった行動のアクションスコアを0に設定する。このような処理を繰り返すことにより、状態Sとアクションスコアとで構成される図57に示すような行動履歴データが得られる。

【0186】

行動履歴データが得られると、情報処理装置10は、図58に示した処理の流れに沿って思考ルーチンを構築する。図58に示すように、行動履歴データを取得すると(S301)、情報処理装置10は、取得した行動履歴データを利用した機械学習により思考ルーチンを構築する(S302)。次いで、情報処理装置10は、必要に応じて追加の行動履歴データを取得する(S303)。次いで、情報処理装置10は、追加した行動履歴データと元の行動履歴データとを統合する(S304)。次いで、情報処理装置10は、終了条件を満たしたか否かを判定する(S305)。

【0187】

10

20

30

40

50

例えば、ユーザによる終了操作が与えられた場合や、ランダムに行動するプレイヤーに対する勝率が所定の閾値を越えた場合などに、情報処理装置10は、終了条件を満たしたと判定する。終了条件を満たしていない場合、情報処理装置10は、処理をステップS302に進める。一方、終了条件を満たした場合、情報処理装置10は、思考ルーチンの構築に係る一連の処理を終了する。

【0188】

このようにして自動構築された思考ルーチンを用いてプレイヤーを行動させた結果、ランダムに行動するプレイヤーに対する勝率について、図59に示すような結果が得られた。図59に示すように、15世代(図58のステップS302~S304の繰り返し回数が15)で思考ルーチンを利用して行動するプレイヤーの勝率が100%近くに達した。なお、行動の選択は、最もアクションスコアの高い行動を選択する方法で行われている。但し、この例では、行動を選択する際に、各アクションスコアに僅かの乱数を加えてから行動を選択するようにしている。

10

【0189】

また、先に説明した複数エージェントの同時学習を適用し、2人のプレイヤーの行動を同時に学習して思考ルーチンを構築してみた。複数エージェントの同時学習を適用すると、ランダムでない動きをするプレイヤーに対して勝とうとする思考ルーチンが自動構築されるため、より賢くプレイヤーを行動させる思考ルーチンが構築される。なお、互いに思考ルーチンを用いて行動する2人のプレイヤーを対戦させた結果を図60に示した。図60に示すように、学習世代によりプレイヤー1が大きく勝ち越す場合もあるが、プレイヤー2が大きく勝ち越す場合もある。

20

【0190】

また、ある学習世代において実験的に1000試合のゲームを行った結果、図61に示すように、プレイヤー1が大きく勝ち越す結果(対戦勝率)が得られた。但し、ランダムに行動するプレイヤーを相手にした場合(ランダム相手)、プレイヤー1もプレイヤー2も相手に対して9割以上の高い勝率を得た。つまり、思考ルーチンを利用して行動するプレイヤーは、十分に賢く行動しているのである。このように、複数エージェントの同時学習を適用すると、相手に勝とうとして思考ルーチンを強化しているうちに、ランダムに行動する相手に対しても高い確率で勝てる汎用的なアルゴリズムが得られる。

30

【0191】

ところで、これまでは状態Sとして、自分の座標、相手の座標、時刻を表現した3次元マトリクスを用いていたが、この3次元マトリクスに代えてゲーム画面の画像情報をそのまま用いる方法も考えられる。例えば、状態Sとして、図62に示すようなゲーム画面の輝度画像を用いることができる。つまり、状態Sは、行動を決定するために有用な情報が含まれてさえいれば何でもよいのである。この考えに基づくと、本実施形態に係る技術が様々なゲームやタスクに関する思考ルーチンの自動構築方法に応用できることが容易に想像できるであろう。

【0192】

以上、「対戦ゲーム」への応用について説明した。

【0193】

[3-3:「五目並べ」への応用]

次に、図63及び図64を参照しながら、本実施形態に係る技術を「五目並べ」へと応用する方法について説明する。「五目並べ」の主なルールは、(1)交互に手を打つ、(2)縦横斜めに先に5つの石を並べた方が勝ち、の2点である。また、「五目並べ」において、状態Sは盤面であり、行動aは各プレイヤーが打つ手である。

40

【0194】

「五目並べ」への応用方法は、基本的に「三目並べ」への応用方法と同じである。つまり、状態S及び行動aは、図63に示すように、2次元マトリクスで表現される。また、最初に用いる行動履歴データは、完全にランダムに石を配置するInitialを用いて取得される。そして、最終的に勝ちに至った全ての(S, a)に報酬1が設定され、負けに

50

至った全ての (S , a) に報酬 0 が設定される。情報処理装置 10 は、この行動履歴データを用いて思考ルーチンを構築する。また、情報処理装置 10 は、思考ルーチンを用いて対局し、その結果を統合した行動履歴データを用いて思考ルーチンを構築する。これらの処理を繰り返すことにより、賢い行動を選択する思考ルーチンが構築される。

【 0 1 9 5 】

また、行動を選択する際、情報処理装置 10 は、「三目並べ」の場合と同様に全ての行動の可能性について (石を置く全ての点について石を置いたとして) 推定報酬を求め、最も推定報酬の高くなる点に石を置く。もちろん、情報処理装置 10 は、数手先を読んで石を置く位置を選択するように構成されていてもよい。なお、「五目並べ」は、「三目並べ」に比べて盤面の組み合わせ数が膨大である。そのため、ランダムに石を置くプレイヤーは、見当違いの手を打ちがちであるために非常に弱い。

10

【 0 1 9 6 】

従って、ランダムに石を置くプレイヤーを相手に学習を行っても非常に弱い相手に勝つための思考ルーチンができあがるだけで、賢い思考ルーチンはなかなか得られない。そこで、対戦ゲームと同様に、複数エージェントの同時学習を適用し、自分と相手とを同じ環境で学習させる手法を用いる方が好ましい。このような構成にすることで、比較的高性能な思考ルーチンを自動構築することが可能になる。互いに思考ルーチンを用いて行動するプレイヤーによる対局結果を図 6 4 に示した。

【 0 1 9 7 】

以上、「五目並べ」への応用について説明した。

20

【 0 1 9 8 】

[3 - 4 : 「ポーカー」への応用]

次に、図 6 5 ~ 図 6 7 を参照しながら、本実施形態に係る技術を「ポーカー」へと応用する方法について説明する。「ポーカー」の主なルールは、図 6 5 に示すように、(1) 5 枚のカードを配る、(2) 捨てるカードを選択する、(3) 役の強い方が勝ち、の 3 点である。ここでは、カードが配られたときに、捨てるカードを決める思考ルーチンの構築方法について考える。

【 0 1 9 9 】

図 6 6 に示すように、状態 S 及び行動 a は、記号列で表現される。例えば、ハートのエースを " H A "、クラブの 2 を " C 2 "、ダイヤの K を " D K " などと表現する。図 6 6 の場合、状態 S は、記号列 " S J C J C 0 D 9 D 7 " で表現される。また、ダイヤの 9 及びダイヤの 7 を捨てた場合、行動 a は、記号列 " D 9 D 7 " で表現される。また、ゲームに勝った場合には報酬 " 1 " が与えられ、負けた場合には報酬 " 0 " が与えられる。このような表現を用いると、例えば、図 6 7 に示すような行動履歴データが得られる。

30

【 0 2 0 0 】

最初に行動履歴データを取得する Innate としては、例えば、完全にランダムに 5 枚のカードそれぞれを捨てるかどうか決定するものを利用する。また、情報処理装置 10 は、勝った行動履歴データには報酬 " 1 " を設定し、負けた行動履歴データには報酬 " 0 " を設定する。そして、情報処理装置 10 は、蓄積された行動履歴データを用いて思考ルーチンを構築する。このとき、行動を選択した結果、どのような役が揃ったか、相手の役がどのようなものだったかなどの情報は利用されない。つまり、純粹に勝ち負けだけを考慮して思考ルーチンが構築される。但し、自分が強い役を揃えるのに有利なカードの切り方を選択した行動履歴データほど報酬が 1 になる確率は高くなる傾向にある。

40

【 0 2 0 1 】

さて、行動を選択する際、配られた 5 枚のカードそれぞれについて、カードを切る、カードを切らない、の選択肢が与えられる。そのため、行動の組み合わせは、2 の 5 乗 = 3 2 通り存在する。従って、思考ルーチンは、報酬推定機を利用し、3 2 通りの (S , a) について推定報酬を算出し、最も推定報酬の高い行動を選択する。

【 0 2 0 2 】

以上、「ポーカー」への応用について説明した。

50

【0203】

[3-5:「ロールプレイングゲーム」への応用]

次に、図68～図79を参照しながら、本実施形態に係る技術を「ロールプレイングゲーム」へと応用する方法について説明する。ここでは、「ロールプレイングゲーム」の戦闘シーンにおいてプレイヤーに代わってキャラクタを賢く自動操作する思考ルーチンの自動構築方法について考える。なお、ここで考える「ロールプレイングゲーム」のルールは、図68に示した通りである。また、図68に示すように、状態Sはプレイヤーに提供される情報であり、行動aはキャラクタを操作するコマンドである。

【0204】

戦闘シーンの環境は、図69に示した通りである。まず、戦闘に勝つと生存者で経験値が山分けされる。さらに、経験値が貯まるとレベルアップする。また、レベルアップすると、キャラクタの職業に応じてステータスの値がアップしたり、魔法を覚えたりする。また、戦闘に5回連続で勝つと敵のレベルが1つアップすると共に、キャラクタのヒットポイントが回復する。また、敵のレベルが31に達するとゲームをクリアしたことになる。

【0205】

なお、戦闘シーンにおいて、キャラクタが持つステータスの1つである“素早さ”の値に応じて各キャラクタが行動をおこせるタイミングが決まる。また、キャラクタがとれる行動は、“攻撃”及び“魔法(魔法を覚えている場合)”である。魔法の種類としては、Heal、Fire、Iceがある。Healは、味方のヒットポイント(HP)を回復する魔法である。Fireは、火を用いて敵を攻撃する魔法である。Iceは、氷を用いて敵を攻撃する魔法である。また、魔法をかけるターゲットは、単体又は全体のいずれかを選択可能である。但し、全体を選択した場合には魔法の効果が半減する。また、使える魔法の種類やレベルは、キャラクタのレベルに応じて変わる。さらに、同じ魔法でもレベルの高い魔法ほどマジックポイント(MP)を多く消費するが、効果は高い。

【0206】

キャラクタの職業及び職業毎のステータスは、図70に示した通りである。ステータス上昇率は、キャラクタのレベルが1つアップする度にステータスがアップする割合を示している。また、魔法を覚えるLvは、記載された値のレベルに達した場合にキャラクタが魔法を覚えるレベルを示している。但し、空欄に対応する魔法は覚えられない。また、0と記載されている箇所は、最初から魔法を覚えていることを示している。なお、味方のパーティは、上側4種類の職業を持つキャラクタにより構成される。一方、敵のパーティは、下側4種類の職業を持つキャラクタから選択されたキャラクタにより構成される。

【0207】

状態Sとして利用される味方側の情報は、図71に示した通りである。例えば、生存する味方のレベル、職業、HP、最大HP、MP、最大MP、攻撃力、防御力、素早さなどが状態Sとして利用される。なお、職業の欄は、当てはまる職業の欄に1、それ以外の欄に0が記入される。また、その他の欄には現状の値が記入される。一方、状態Sとして利用される敵側の情報は、図72に示した通りである。例えば、生存する敵のレベル、職業、累積ダメージなどが状態Sとして利用される。なお、累積ダメージは、それまでに与えたダメージの合計値を示している。

【0208】

また、行動aとして利用される味方側の情報は、図73に示した通りである。例えば、行動者の欄には、これから行動を行うキャラクタの場合に1、それ以外の場合に0が記入される。また、行動対象の欄には、行動の対象となるキャラクタの場合に1、それ以外の場合に0が記入される。例えば、回復魔法を受けるキャラクタに対応する行動対象の欄には1が記入される。また、アクションの種類欄には、行う行動の欄に1、行わない行動の欄に0が記入される。一方、行動aとして利用される敵側の情報は、図74に示した通りである。図74に示すように、敵側の情報としては行動対象の情報が利用される。

【0209】

さて、これまで説明してきた応用例と同様、情報処理装置10は、まず、行動履歴デー

10

20

30

40

50

タを取得する。このとき、情報処理装置10は、行動の種類毎に選択確率に重みを付けた上で、ランダムに行動を選択するInnateを用いて行動履歴データを取得する。例えば、情報処理装置10は、魔法よりも攻撃を選択する確率を高く設定したInnateを用いてキャラクタを行動させる。また、図75に示すように、情報処理装置10は、味方がやられた場合には報酬“-5”を行動履歴データに設定し、敵を倒した場合には報酬“1”を行動履歴データに設定する。その結果、図76のA図に示すような行動履歴データが得られる。但し、味方や敵がやられる過程の評価も考慮するため、情報処理装置10は、図76のB図に示すように、直線的に報酬の値をDecayさせる。

【0210】

情報処理装置10は、上記のようにして取得された行動履歴データを用いて思考ルーチンを構築する。このとき、情報処理装置10は、時刻tにおける状態S及び行動aから時刻t+1における状態S'を推定する予測機を構築する。また、情報処理装置10は、時刻t+1における状態S'から推定報酬を算出する報酬推定機を構築する。そして、情報処理装置10は、図77に示すように、現在の状態Sにおいてキャラクタがとりうる行動毎に、予測機を用いて次の状態S'を予測する。さらに、情報処理装置10は、予測した状態S'を報酬推定機に入力して推定報酬yを算出する。推定報酬yを算出した情報処理装置10は、推定報酬yが最大となる行動aを選択する。

【0211】

図77の例では、行動a = “敵全体にFire”に対応する推定報酬yが最大となっている。そのため、この例においては、最適な行動として、行動a = “敵全体にFire”が選択される。但し、思考ルーチンは、図78に示すように、推定報酬が高く、推定誤差が大きく、特徴量空間における密度が疎な特徴量座標に対応する行動を選択するように構成されていてもよい。つまり、先に説明した効率的な推定報酬機の構築方法で紹介した報酬スコア、未知スコア、誤差スコアに基づいて思考ルーチンが構築されていてもよい。

【0212】

なお、報酬スコアは、とりうる全ての行動について報酬推定機を用いて推定報酬を求め、推定報酬の低い方から順に1、2、3、...と、推定報酬が高くなるほど大きくなるようにスコアを与えることで得られる。また、未知スコアは、図25などに示した方法を用いて、全ての行動について特徴量座標の周辺密度を求め、密度が高い方から順に1、2、3、...と、密度が低くなるほど大きくなるようにスコアを与えることで得られる。

【0213】

また、誤差スコアを求める場合、情報処理装置10は、まず、既存の行動履歴データの全てについて、推定報酬yの値を実際の報酬rと比較し、その誤差を求める。次いで、情報処理装置10は、平均値よりも誤差の大きい行動履歴データに対応する特徴量座標を特徴量空間にプロットする。次いで、情報処理装置10は、プロットした特徴量座標の密度分布を求める。最後に、情報処理装置10は、全ての行動履歴データに対応する特徴量座標について密度を求め、密度が低い方から順に1、2、3、...と、密度が高くなるほど大きくなるようにスコアを与える。

【0214】

例えば、報酬スコアを s_1 、未知スコアを s_2 、誤差スコアを s_3 と表記した場合、情報処理装置10は、行動を選択する際に、 $s_1 * w_1 + s_2 * w_2 + s_3 * w_3$ （但し、 $w_1 \sim w_3$ は所定の重み）の値を算出し、この値が最も大きくなる行動を選択する。このようにして行動を選択することにより、報酬が高く、推定誤差が大きく、特徴量空間における特徴量座標の密度が疎な行動を選択することが可能になる。

【0215】

ここで、図79を参照しながら、効率的な推定報酬機の構築方法を適用した場合の効果について述べる。図79のグラフは、最も高い推定報酬が得られる行動を選択した場合（最適戦略）と、効率的な推定報酬機の構築方法を適用した場合（探索行動）とでシナリオクリアまでの1ステップ当たりの平均Rewardを比較したグラフである。図79のグラフから明らかなように、3つのスコアを利用して構築された思考ルーチン（探索行動）

10

20

30

40

50

の方が安定して高い報酬が得られている。この評価結果から、効率的な推定報酬機の構築方法を適用することで、演算負荷を軽減できるばかりか、より高性能な思考ルーチンを構築できることが分かった。

【0216】

なお、「ロールプレイングゲーム」に応用すべく自動構築された思考ルーチンは、以下のような戦略を自身で身に付けていることも分かった。

(A) 集中攻撃：

集中攻撃することで敵の数を素早く減らす。

(B) HP 減少したら回復：

HP が減少した味方のHP を回復して味方が倒されにくくする。

(C) 単体攻撃、全体攻撃の使い分け：

敵の数がある程度多い時は全体攻撃魔法を使う。現在集中攻撃中の敵が残りわずかなダメージで倒せそうなときは、全体攻撃魔法を用いて集中攻撃中の敵を倒しながら、他の敵にもダメージを与える。

(D) 魔法の無駄打ちはしない：

HP の減っていない味方に対して回復魔法を使わない。魔法の効かない敵に対して魔法を使わない。

【0217】

以上、「ロールプレイングゲーム」への応用について説明した。

【0218】

以上説明したように、本実施形態に係る技術を適用すると、人手による調整を介さずに様々な思考ルーチンを自動構築することが可能になる。

【0219】

< 4 : ハードウェア構成例 >

上記の情報処理装置10が有する各構成要素の機能は、例えば、図80に示すハードウェア構成を用いて実現することが可能である。つまり、当該各構成要素の機能は、コンピュータプログラムを用いて図80に示すハードウェアを制御することにより実現される。なお、このハードウェアの形態は任意であり、例えば、パーソナルコンピュータ、携帯電話、PHS、PDA等の携帯情報端末、ゲーム機、又は種々の情報家電がこれに含まれる。但し、上記のPHSは、Personal Handy-phone Systemの略である。また、上記のPDAは、Personal Digital Assistantの略である。

【0220】

図80に示すように、このハードウェアは、主に、CPU902と、ROM904と、RAM906と、ホストバス908と、ブリッジ910と、を有する。さらに、このハードウェアは、外部バス912と、インターフェース914と、入力部916と、出力部918と、記憶部920と、ドライブ922と、接続ポート924と、通信部926と、を有する。但し、上記のCPUは、Central Processing Unitの略である。また、上記のROMは、Read Only Memoryの略である。そして、上記のRAMは、Random Access Memoryの略である。

【0221】

CPU902は、例えば、演算処理装置又は制御装置として機能し、ROM904、RAM906、記憶部920、又はリムーバブル記録媒体928に記録された各種プログラムに基づいて各構成要素の動作全般又はその一部を制御する。ROM904は、CPU902に読み込まれるプログラムや演算に用いるデータ等を格納する手段である。RAM906には、例えば、CPU902に読み込まれるプログラムや、そのプログラムを実行する際に適宜変化する各種パラメータ等が一時的又は永続的に格納される。

【0222】

これらの構成要素は、例えば、高速なデータ伝送が可能なホストバス908を介して相互に接続される。一方、ホストバス908は、例えば、ブリッジ910を介して比較的デ

10

20

30

40

50

ータ伝送速度が低速な外部バス912に接続される。また、入力部916としては、例えば、マウス、キーボード、タッチパネル、ボタン、スイッチ、及びレバー等が用いられる。さらに、入力部916としては、赤外線やその他の電波を利用して制御信号を送信することが可能なりモートコントローラ（以下、リモコン）が用いられることもある。

【0223】

出力部918としては、例えば、CRT、LCD、PDP、又はELD等のディスプレイ装置、スピーカ、ヘッドホン等のオーディオ出力装置、プリンタ、携帯電話、又はファクシミリ等、取得した情報を利用者に対して視覚的又は聴覚的に通知することが可能な装置である。但し、上記のCRTは、Cathode Ray Tubeの略である。また、上記のLCDは、Liquid Crystal Displayの略である。そして、上記のPDPは、Plasma Display Panelの略である。さらに、上記のELDは、Electro-Luminescence Displayの略である。

10

【0224】

記憶部920は、各種のデータを格納するための装置である。記憶部920としては、例えば、ハードディスクドライブ（HDD）等の磁気記憶デバイス、半導体記憶デバイス、光記憶デバイス、又は光磁気記憶デバイス等が用いられる。但し、上記のHDDは、Hard Disk Driveの略である。

【0225】

ドライブ922は、例えば、磁気ディスク、光ディスク、光磁気ディスク、又は半導体メモリ等のリムーバブル記録媒体928に記録された情報を読み出し、又はリムーバブル記録媒体928に情報を書き込む装置である。リムーバブル記録媒体928は、例えば、DVDメディア、Blu-rayメディア、HD DVDメディア、各種の半導体記憶メディア等である。もちろん、リムーバブル記録媒体928は、例えば、非接触型ICチップを搭載したICカード、又は電子機器等であってもよい。但し、上記のICは、Integrated Circuitの略である。

20

【0226】

接続ポート924は、例えば、USBポート、IEEE1394ポート、SCSI、RS-232Cポート、又は光オーディオ端子等のような外部接続機器930を接続するためのポートである。外部接続機器930は、例えば、プリンタ、携帯音楽プレーヤ、デジタルカメラ、デジタルビデオカメラ、又はICレコーダ等である。但し、上記のUSBは、Universal Serial Busの略である。また、上記のSCSIは、Small Computer System Interfaceの略である。

30

【0227】

通信部926は、ネットワーク932に接続するための通信デバイスであり、例えば、有線又は無線LAN、Bluetooth（登録商標）、又はWUSB用の通信カード、光通信用のルータ、ADSL用のルータ、又は各種通信用のモデム等である。また、通信部926に接続されるネットワーク932は、有線又は無線により接続されたネットワークにより構成され、例えば、インターネット、家庭内LAN、赤外線通信、可視光通信、放送、又は衛星通信等である。但し、上記のLANは、Local Area Networkの略である。また、上記のWUSBは、Wireless USBの略である。そして、上記のADSLは、Asymmetric Digital Subscriber Lineの略である。

40

【0228】

以上、ハードウェア構成例について説明した。

【0229】

<5:まとめ>

最後に、本実施形態の技術的思想について簡単に纏める。以下に記載する技術的思想は、例えば、PC、携帯電話、携帯ゲーム機、携帯情報端末、情報家電、カーナビゲーションシステム等、種々の情報処理装置に対して適用することができる。

【0230】

50

上記の情報処理装置の機能構成は以下のように表現することができる。例えば、下記（１）に記載の情報処理装置は、行動履歴データを用いて報酬推定機を自動構築することができる。この報酬推定機を利用すると、エージェントがおかれた状態に応じて、その状態でエージェントがとりうる行動毎に、行動を行ったエージェントが得る報酬を推定することができる。そのため、高い報酬を得ると推定される行動をエージェントがとるように制御することで、賢く行動するエージェントの動きを実現することが可能になる。言い換えると、下記（１）に記載の情報処理装置は、賢く行動するエージェントの動きを実現することが可能な思考ルーチンを自動構築することができる。

【 0 2 3 1 】

（ 1 ）

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成する報酬推定機生成部を備え、

前記報酬推定機生成部は、

複数の処理関数を組み合わせることで複数の基底関数を生成する基底関数生成部と、

前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出部と、

前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰 / 判別学習により算出する推定関数算出部と、

を含み、

前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、

情報処理装置。

【 0 2 3 2 】

（ 2 ）

エージェントがおかれた現在の状態を表す状態データと、当該エージェントが次にとりうる行動を表す行動データとを前記報酬推定機に入力し、当該行動をとった結果として当該エージェントが得る報酬値を推定する報酬値推定部と、

前記現在の状態において前記エージェントが次にとりうる行動のうち、前記報酬値推定部により推定された報酬値が最も高い値となる行動を選択する行動選択部と、

をさらに備える、

上記（ 1 ）に記載の情報処理装置。

【 0 2 3 3 】

（ 3 ）

前記行動選択部による選択結果に基づいてエージェントを行動させる行動制御部と、

前記エージェントの行動に伴って更新される状態データ及び行動データを蓄積し、当該行動の結果としてエージェントが得た報酬を表す報酬値と、蓄積した状態データ及び行動データとを対応付けて前記行動履歴データに追加する履歴データ追加部と、

をさらに備える、

上記（ 2 ）に記載の情報処理装置。

【 0 2 3 4 】

（ 4 ）

前記状態データ、前記行動データ、及び前記報酬値の組が前記行動履歴データに追加された場合、前記特徴量ベクトル算出部は、前記行動履歴データに含まれる全ての状態データ及び行動データについて特徴量ベクトルを算出し、

前記情報処理装置は、特徴領空間において前記特徴量ベクトルにより示される座標点の分布が所定の分布に近づくように前記行動履歴データに含まれる前記状態データ、前記行動データ、及び前記報酬値の組を間引く分布調整部をさらに備える、

上記（ 3 ）に記載の情報処理装置。

【 0 2 3 5 】

10

20

30

40

50

(5)

前記状態データ、前記行動データ、及び前記報酬値の組が前記行動履歴データに追加された場合、前記特徴量ベクトル算出部は、前記行動履歴データに含まれる全ての状態データ及び行動データについて特徴量ベクトルを算出し、

前記情報処理装置は、特徴領空間において前記特徴量ベクトルにより示される座標点の分布が所定の分布に近づくように前記行動履歴データに含まれる前記状態データ、前記行動データ、及び前記報酬値の組のそれぞれに重みを設定する分布調整部をさらに備える、

上記(3)に記載の情報処理装置。

【 0 2 3 6 】

(6)

前記分布調整部は、間引き後に残った前記状態データ、前記行動データ、及び前記報酬値の組について、特徴領空間において前記特徴量ベクトルにより示される座標点の分布が所定の分布に近づくように前記行動履歴データに含まれる前記状態データ、前記行動データ、及び前記報酬値の組のそれぞれに重みを設定する、

上記(4)に記載の情報処理装置。

【 0 2 3 7 】

(7)

前記行動履歴データを学習データとして用い、現在の時刻においてエージェントがおかれた状態を表す状態データ及び現在の時刻においてエージェントがとる行動を表す行動データから次の時刻におけるエージェントの状態を表す状態データを予測する予測機を機械学習により生成する予測機生成部をさらに備え、

前記報酬値推定部は、

現在の時刻における状態データ及び行動データを前記予測機に入力して次の時刻におけるエージェントの状態を表す状態データを予測し、

前記次の時刻におけるエージェントの状態を表す状態データと、当該状態においてエージェントがとりうる行動を表す行動データとを前記報酬推定機に入力して、当該行動をとった結果として当該エージェントが得る報酬値を推定する、

上記(2) ~ (6) のいずれか 1 項に記載の情報処理装置。

【 0 2 3 8 】

(8)

前記行動履歴データを学習データとして用い、現在の時刻においてエージェントがおかれた状態を表す状態データ及び現在の時刻においてエージェントがとる行動を表す行動データから次の時刻におけるエージェントの状態を表す状態データを予測する予測機を機械学習により生成する予測機生成部をさらに備え、

前記報酬値推定部は、現在の時刻を時刻 t_0 とした場合に、

時刻 t_0 における状態データ及び行動データを前記予測機に入力して次の時刻 t_1 におけるエージェントの状態を表す状態データを予測する処理を実行し、

$k = 1 \sim n - 1$ ($n \geq 2$) について、時刻 t_k における状態データ及び時刻 t_k においてエージェントがとりうる行動を表す行動データを前記予測機に入力して時刻 t_{k+1} におけるエージェントの状態を表す状態データを予測する処理を逐次実行し、

予測した時刻 t_n におけるエージェントの状態を表す状態データと、当該状態においてエージェントがとりうる行動を表す行動データとを前記報酬推定機に入力して、当該行動をとった結果として当該エージェントが得る報酬値を推定する、

上記(2) ~ (6) のいずれか 1 項に記載の情報処理装置。

【 0 2 3 9 】

(9)

前記報酬推定機生成部は、複数のエージェントの状態を表す状態データと、当該状態において各エージェントがとった行動を表す行動データと、当該行動の結果として各エージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成す

10

20

30

40

50

る、

上記(1)～(8)のいずれか1項に記載の情報処理装置。

【0240】

(10)

前記基底関数生成部は、遺伝的アルゴリズムに基づいて前記基底関数を更新し、

前記特徴量ベクトル算出部は、前記基底関数が更新された場合に、更新後の前記基底関数に前記状態データ及び前記行動データを入力して特徴量ベクトルを算出し、

前記推定関数算出部は、前記更新後の基底関数を用いて算出された特徴量ベクトルの入力に応じて前記報酬値を推定する推定関数を算出する、

上記(1)～(9)のいずれか1項に記載の情報処理装置。

10

【0241】

(11)

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するスコア算出部と、

前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するスコア推定機生成部と、

20

を備え、

前記スコア推定機生成部は、

複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成部と、

前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出部と、

前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習により算出する推定関数算出部と、

を含み、

前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、

30

情報処理装置。

【0242】

(12)

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成するステップを含み、

前記生成するステップは、

複数の処理関数を組み合わせて複数の基底関数を生成するステップと、

前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出するステップと、

40

前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出するステップと、

を含み、

前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、

情報処理方法。

【0243】

(13)

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含

50

む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するステップと、

前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するステップと、

を含み、

前記生成するステップは、

複数の処理関数を組み合わせて複数の基底関数を生成するステップと、

前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出するステップと、

前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習により算出するステップと、

を含み、

前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、情報処理方法。

【0244】

(14)

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データを学習データとして用い、入力された状態データ及び行動データから報酬値を推定する報酬推定機を機械学習により生成する報酬推定機生成機能をコンピュータに実現させるためのプログラムであり、

前記報酬推定機生成機能は、

複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成機能と、

前記行動履歴データに含まれる状態データ及び行動データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出機能と、

前記特徴量ベクトルから前記行動履歴データに含まれる報酬値を推定する推定関数を回帰/判別学習により算出する推定関数算出機能と、

を含み、

前記報酬推定機は、前記複数の基底関数と前記推定関数とにより構成される、プログラム。

【0245】

(15)

エージェントの状態を表す状態データと、当該状態においてエージェントがとった行動を表す行動データと、当該行動の結果としてエージェントが得た報酬を表す報酬値とを含む行動履歴データに基づき、高い報酬を得たエージェントがとった行動及び低い報酬を得たエージェントがとらなかった行動に高いスコアを与え、高い報酬を得たエージェントがとらなかった行動及び低い報酬を得たエージェントがとった行動に低いスコアを与える条件で、各状態データに対応する行動毎のスコアを算出するスコア算出機能と、

前記行動履歴データ及び前記行動毎のスコアを学習データとして用い、入力された状態データから行動毎のスコアを推定するスコア推定機を機械学習により生成するスコア推定機生成機能と、

をコンピュータに実現させるためのプログラムであり、

前記スコア推定機生成機能は、

複数の処理関数を組み合わせて複数の基底関数を生成する基底関数生成機能と、

前記行動履歴データに含まれる状態データを前記複数の基底関数に入力して特徴量ベクトルを算出する特徴量ベクトル算出機能と、

前記特徴量ベクトルから前記行動毎のスコアを推定する推定関数を回帰/判別学習によ

10

20

30

40

50

り算出する推定関数算出機能と、
を含み、

前記スコア推定機は、前記複数の基底関数と前記推定関数とにより構成される、
プログラム。

【0246】

(備考)

上記の報酬推定機構築部12は、報酬推定機生成部の一例である。上記の基底関数リスト生成部121は、基底関数生成部の一例である。上記の特徴量計算部122は、特徴量ベクトル算出部の一例である。上記の推定関数生成部123は、推定関数算出部の一例である。上記の行動選択部14は、報酬値推定部、行動選択部、行動制御部の一例である。上記の行動履歴データ取得部11は、履歴データ追加部の一例である。上記の行動履歴データ統合部124は、分布調整部の一例である。上記の報酬推定機構築部12は、予測機生成部の一例である。上記の報酬推定機構築部12は、スコア算出部、スコア推定機生成部の一例である。

10

【0247】

以上、添付図面を参照しながら本技術に係る好適な実施形態について説明したが、本技術はここで開示した構成例に限定されないことは言うまでもない。当業者であれば、特許請求の範囲に記載された範疇内において、各種の変更例又は修正例に想到し得ることは明らかであり、それらについても当然に本技術の技術的範囲に属するものと了解される。

【符号の説明】

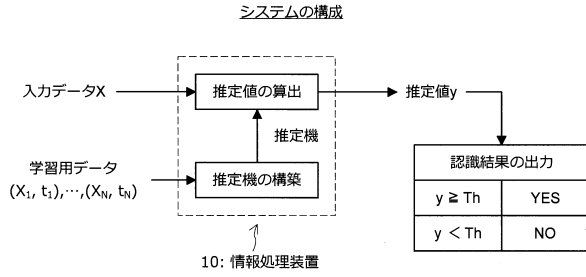
20

【0248】

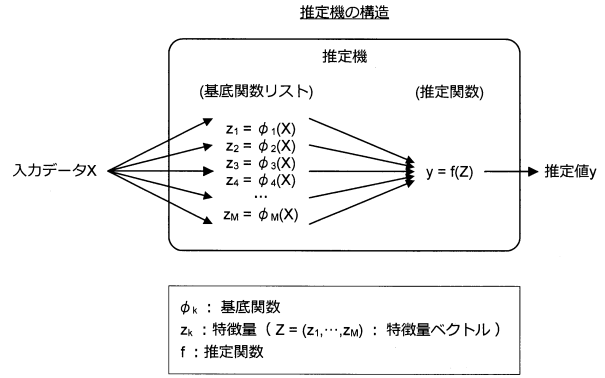
- 10 情報処理装置
- 11 行動履歴データ取得部
- 12 報酬推定機構築部
- 121 基底関数リスト生成部
- 122 特徴量計算部
- 123 推定関数生成部
- 124 行動履歴データ統合部
- 13 入力データ取得部
- 14 行動選択部

30

【図1】



【図3】

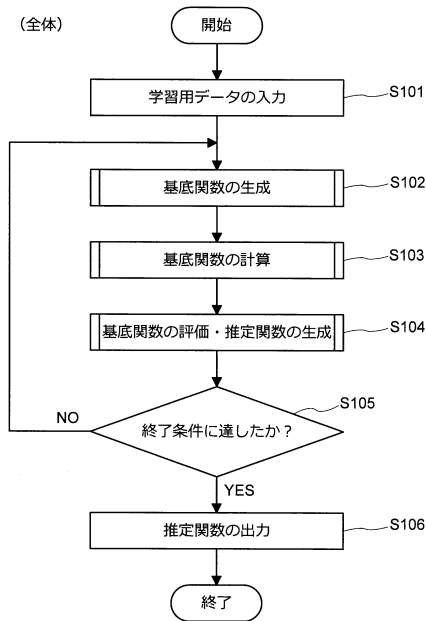


【図2】

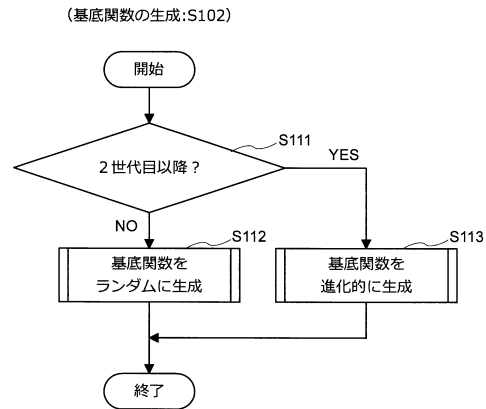
学習用データの例

データ X_k	目的変数 t_k	意味
画像#1	1	画像#1=海
画像#2	1	画像#2=海
⋮	⋮	⋮
画像#N	0	画像#N≠海

【図4】

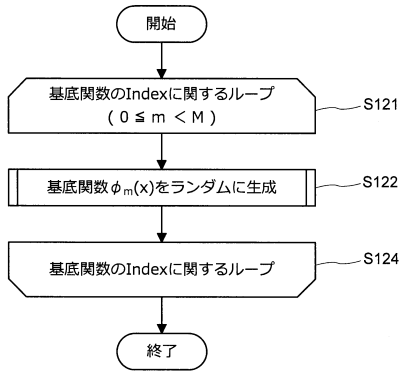


【図5】



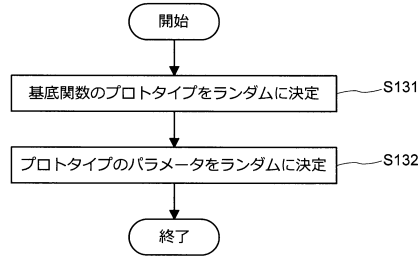
【図6】

(基底関数をランダムに生成:S112)



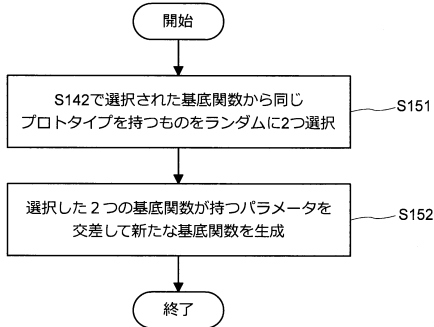
【図7】

(基底関数 φ_m(x)をランダムに生成:S122,S146)



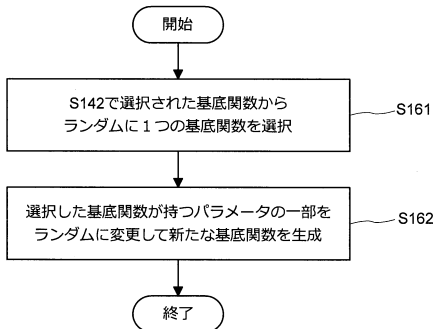
【図9】

(交差:S144)



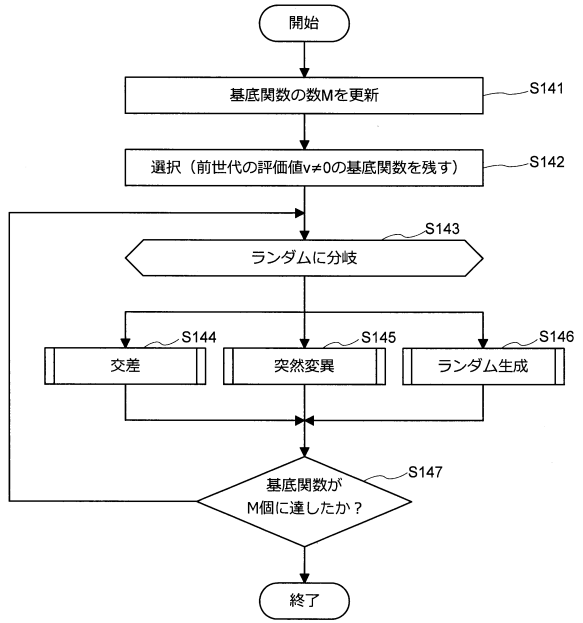
【図10】

(突然変異:S145)



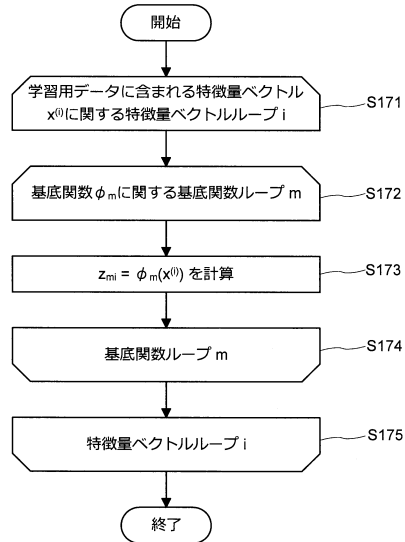
【図8】

(基底関数を進化的に生成:S113)



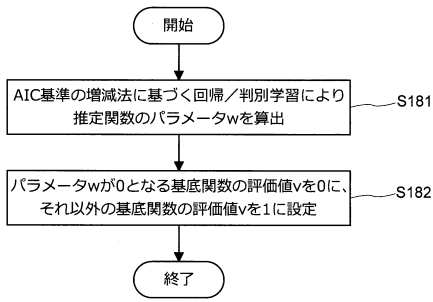
【図11】

(基底関数の計算:S103)

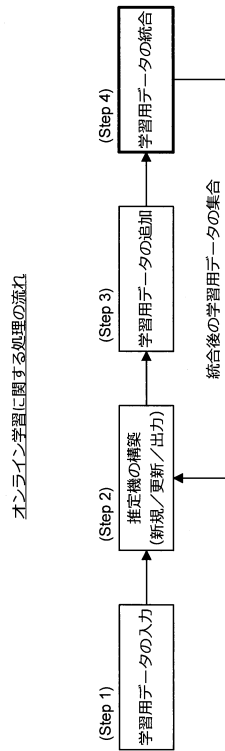


【図12】

(基底関数の評価・推定関数の生成:S104)

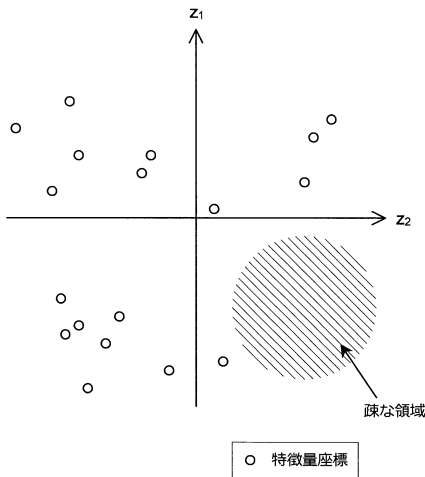


【図13】



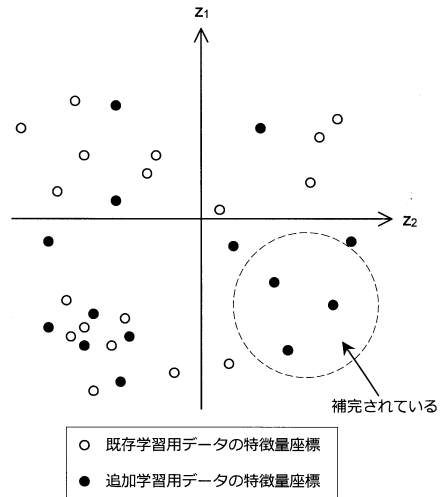
【図14】

特徴量空間における学習用データの分布 (例)

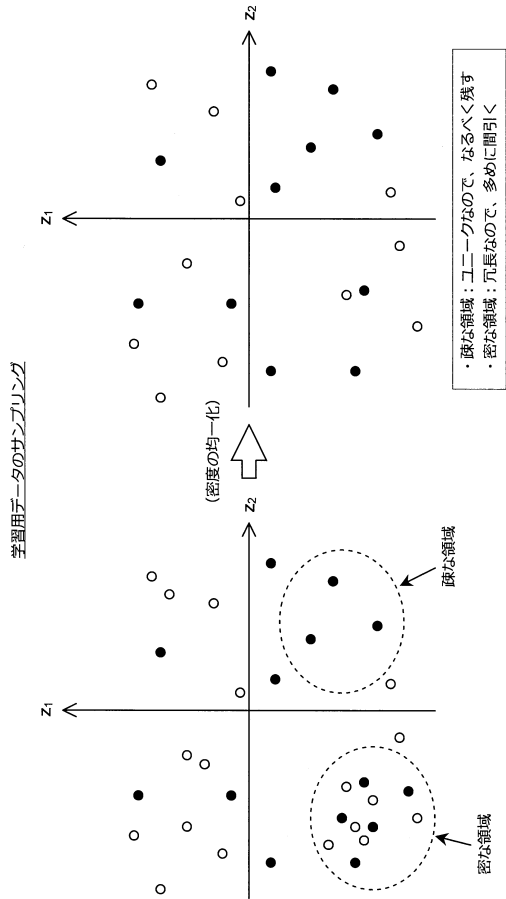


【図15】

特徴量空間における学習用データの分布 (例)



【図16】



【図18】

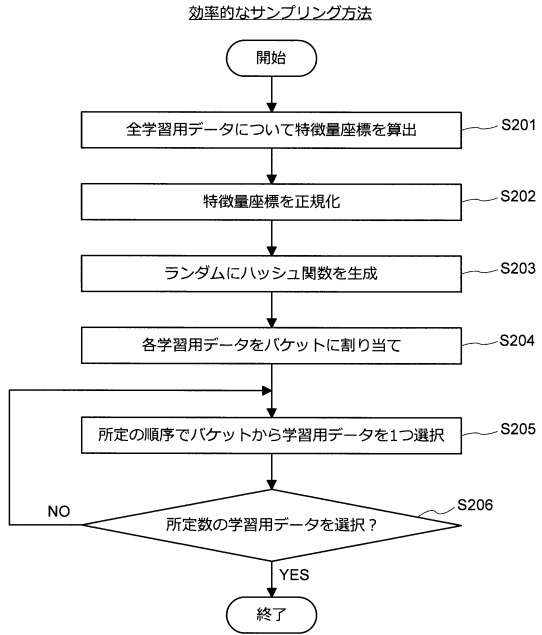
S202: 特徴量座標を正規化

正規化前				
	z_1	z_2	z_3	...
学習用データ#1	12	0.01	-1028	...
学習用データ#2	15	0.03	-650	...
学習用データ#3	11	0.02	-20	...
学習用データ#4	13	0.05	-34	...
⋮	⋮	⋮	⋮	⋮

↓ (分散1、平均0)

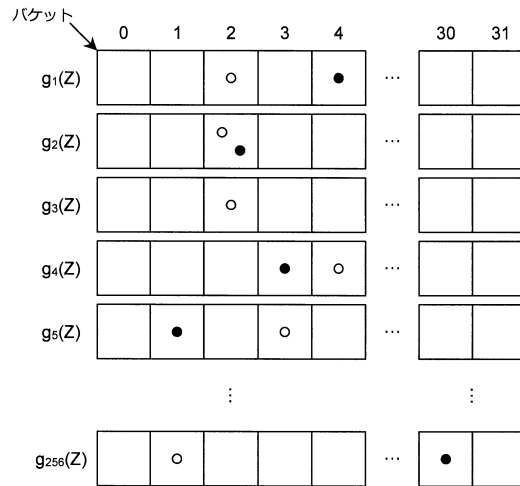
正規化後				
	z_1	z_2	z_3	...
学習用データ#1	-0.43916	-1.02470	-1.20546	...
学習用データ#2	1.31747	0.14639	-0.43964	...
学習用データ#3	-1.02470	-0.43916	0.83673	...
学習用データ#4	0.14639	1.31747	0.80837	...
⋮	⋮	⋮	⋮	⋮

【図17】



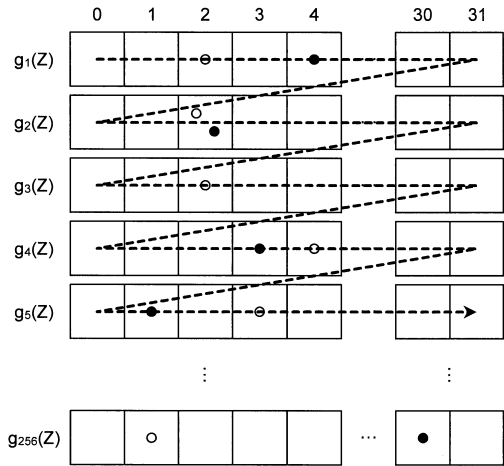
【図19】

S204: 各学習用データをバケットに割り当て



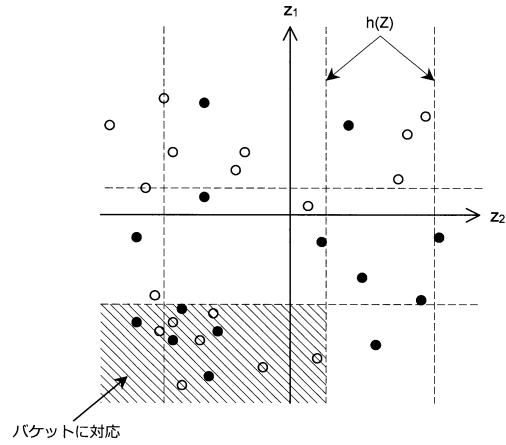
【図20】

S205: 所定の順序でバケットから学習用データを1つ選択



【図22】

特徴量空間とバケットとの対応関係

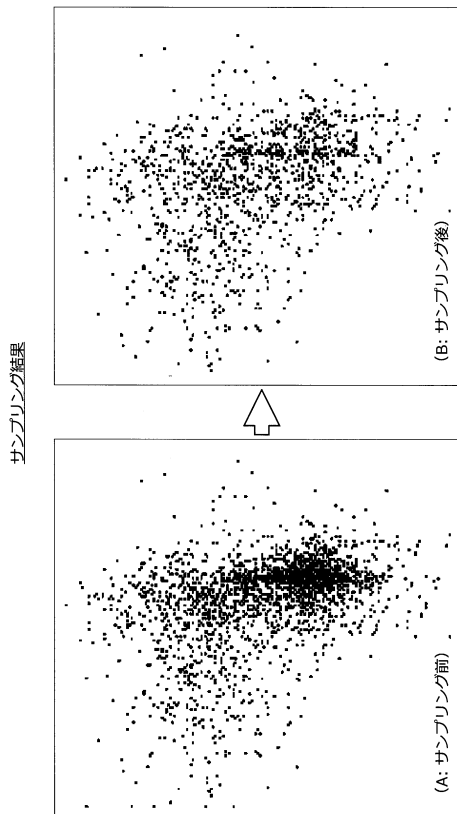


【図21】

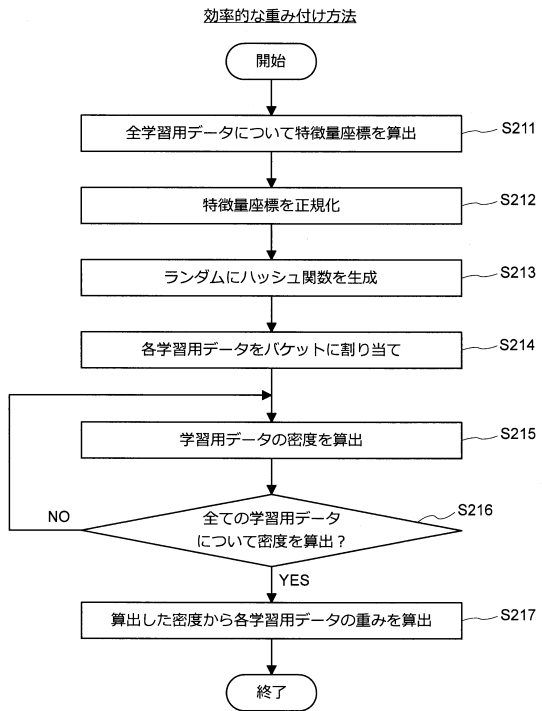
S205: 所定の順序でバケットから学習用データを1つ選択

- <選択ルール>
- ・空のバケットはスキップ
 - ・選択された学習用データを全てのバケットから除く
 - ・1つのバケットに複数の学習用データが割り当てられている場合には、ランダムに1つの学習用データを選択

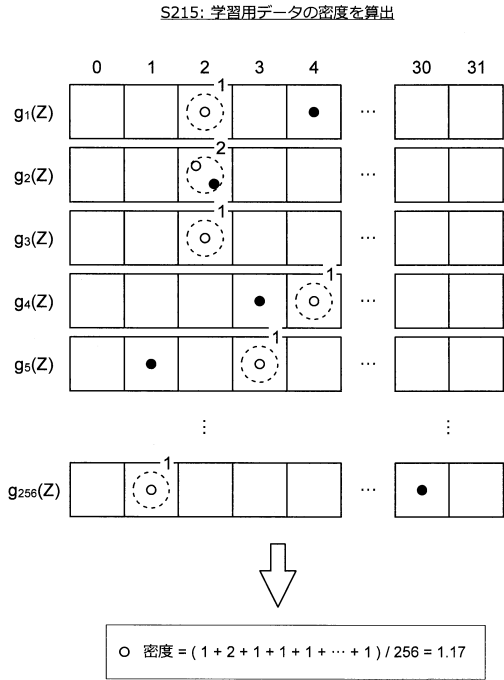
【図23】



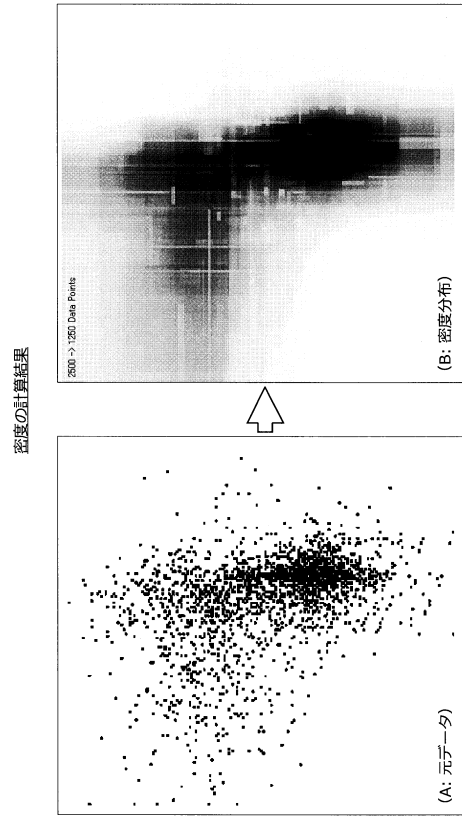
【図24】



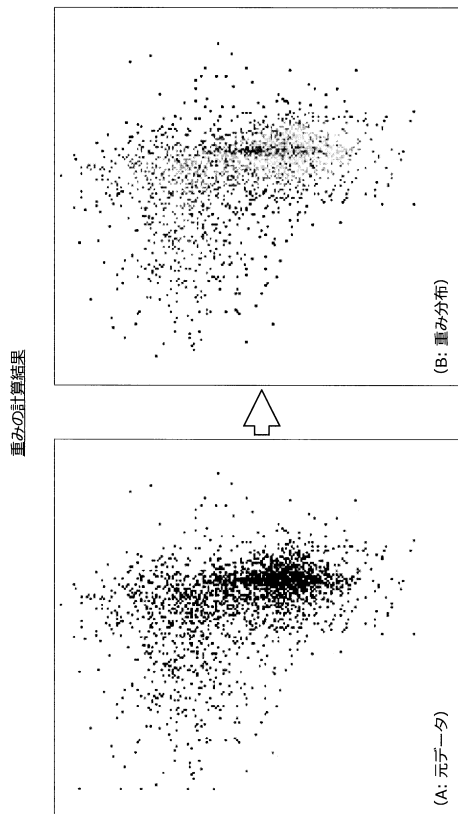
【図 25】



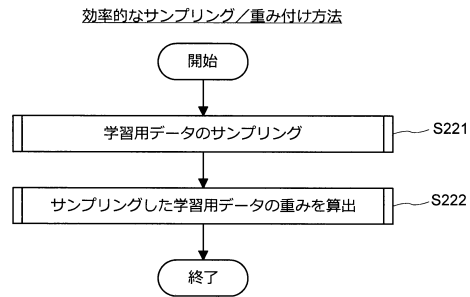
【図 26】



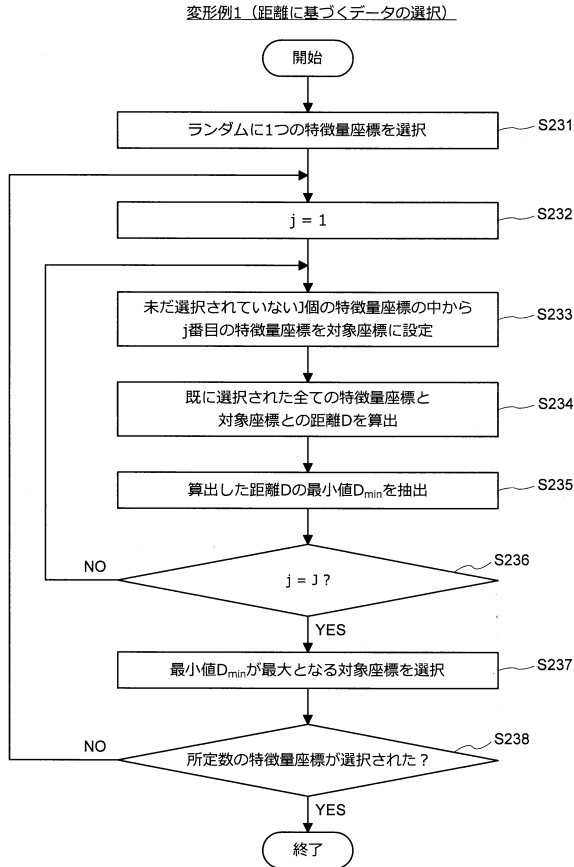
【図 27】



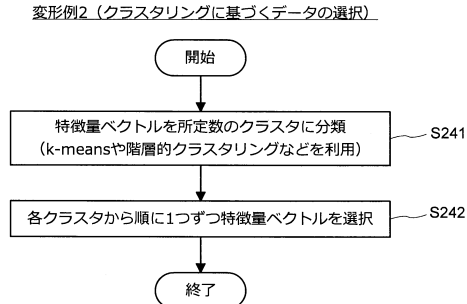
【図 28】



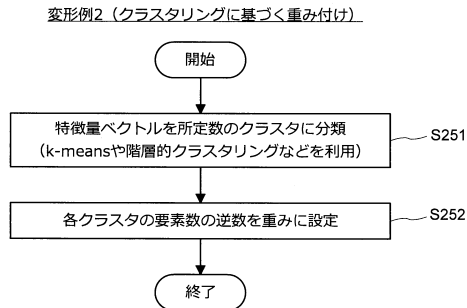
【図29】



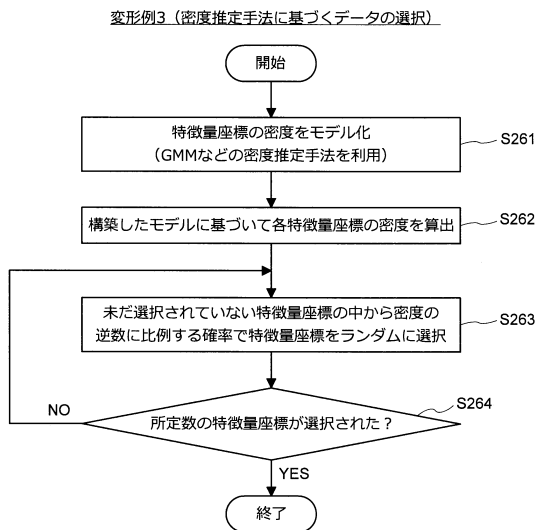
【図30】



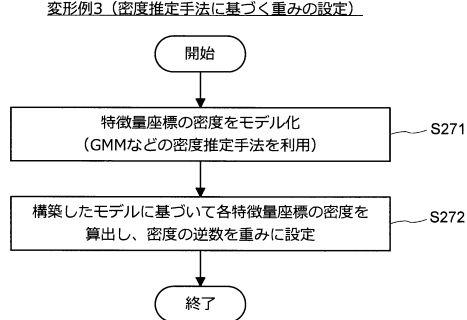
【図31】



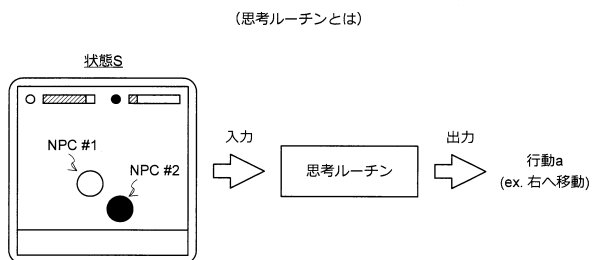
【図32】



【図33】

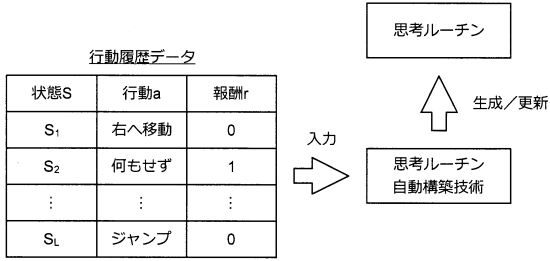


【図34】

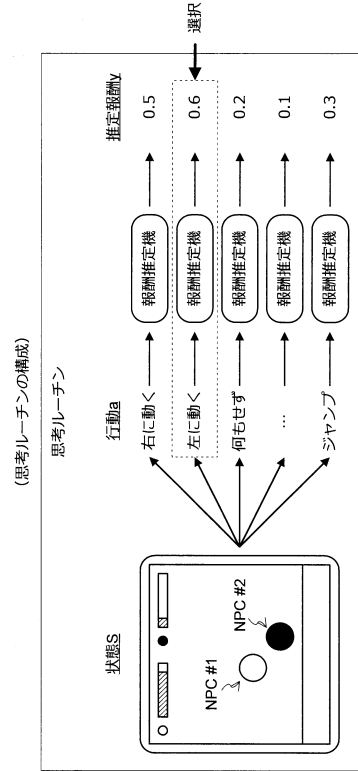


【図 35】

(思考ルーチンの構成)

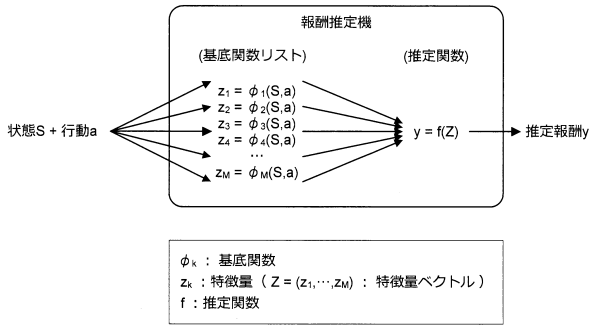


【図 36】



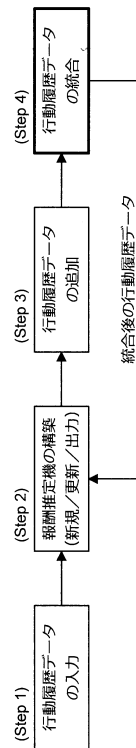
【図 37】

(思考ルーチンの構成; 報酬推定機の構造)

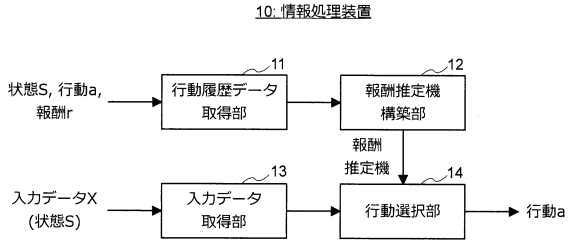


【図 38】

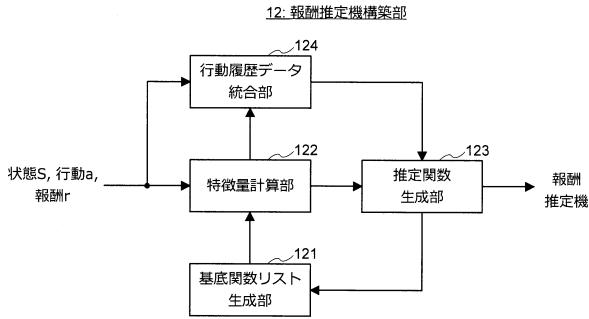
(報酬推定機の構築方法)



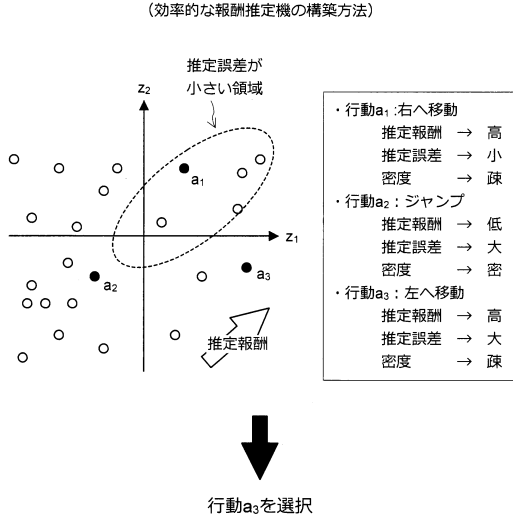
【図39】



【図40】



【図41】



【図42】

(効率的な報酬推定機の構築方法：行動選択時に考慮するスコアの種類)

報酬スコア	推定報酬の高さを表すスコア
未知スコア	特徴量空間における密度の大きさを表すスコア
誤差スコア	推定誤差の大きさを表すスコア

【図43】

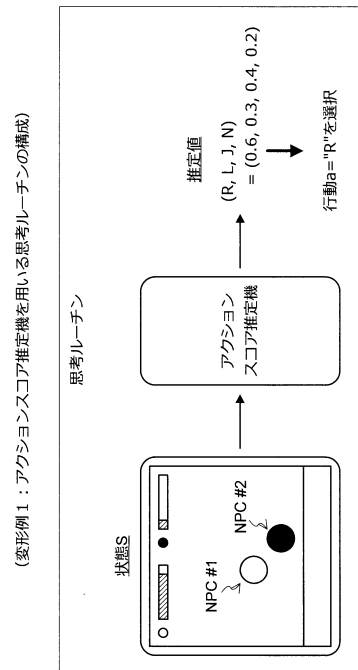
(変形例1：アクションスコア推定機を用いる思考ルーチン)

行動履歴データ

状態S	行動a	報酬r	アクションスコア (R,L,N,J)
S ₁	右へ移動(a="R")	0	(0, 1, 1, 1)
S ₂	左へ移動(a="L")	1	(0, 1, 0, 0)
S ₃	何もせず(a="N")	0	(1, 1, 0, 1)
S ₄	ジャンプ(a="J")	1	(0, 0, 0, 1)

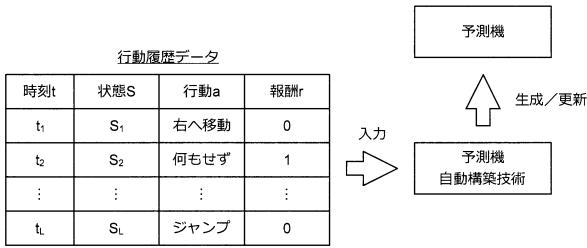


【図44】



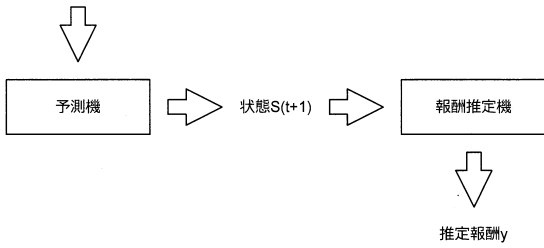
【図45】

(変形例2：予測機を用いた報酬の推定)



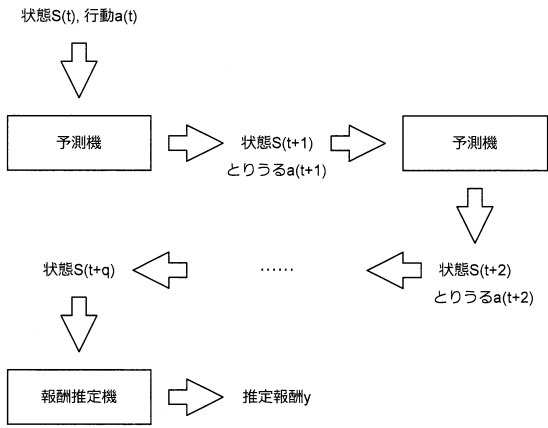
【図46】

状態S(t), 行動a(t) (変形例2：予測機を用いた報酬の推定)



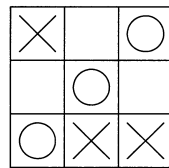
【図47】

(変形例2：予測機を用いた報酬の推定)



【図48】

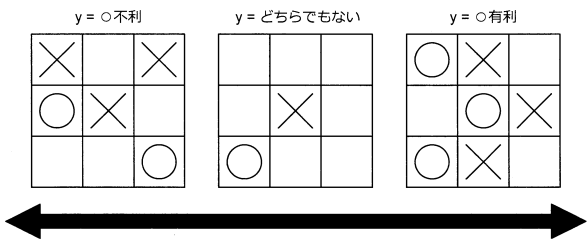
(「三目並べ」への応用)



- <三目並べのルール>
- ・互いに交互に手を打つ
 - ・先に3つ1列に並べた方が勝ち
- 盤面 = S、手 = a

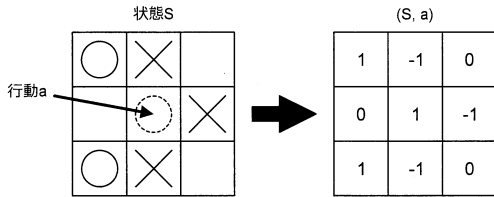
【図49】

(「三目並べ」への応用)



【図51】

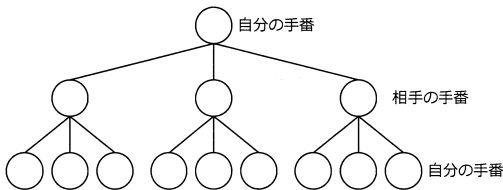
(「三目並べ」への応用)



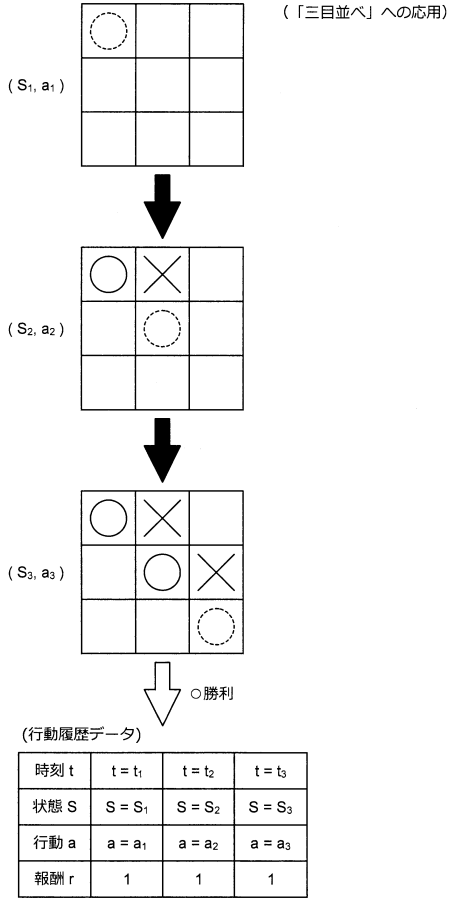
- ・状態S、行動aを3×3のマトリックスで表現
- 状態Sは、自分の手番となった時点の盤面
- 自分の手を反映した盤面が(S, a)
- ・○(自分)を"1"、×(相手)を"-1"、空白を"0"で表現

【図50】

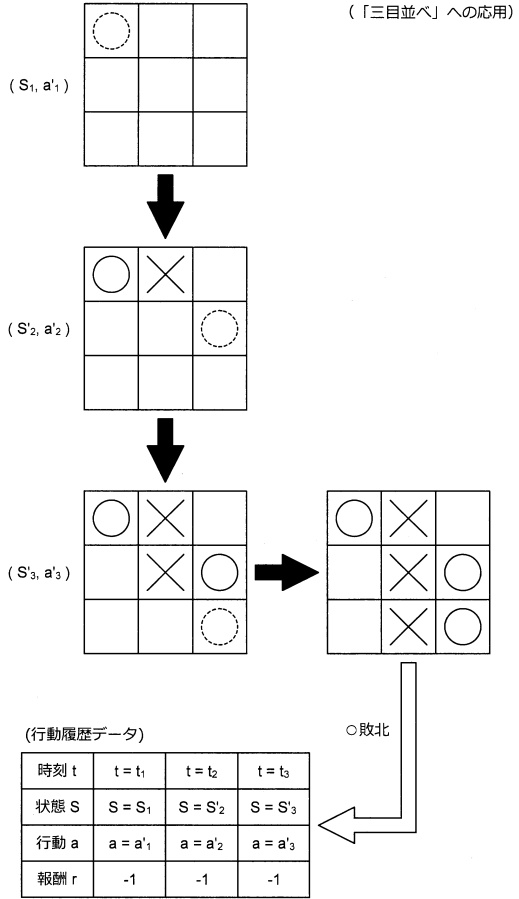
(「三目並べ」への応用)



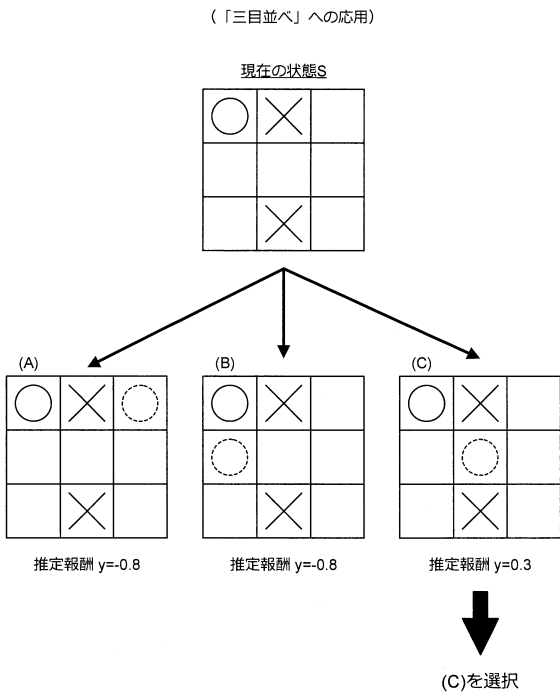
【図52】



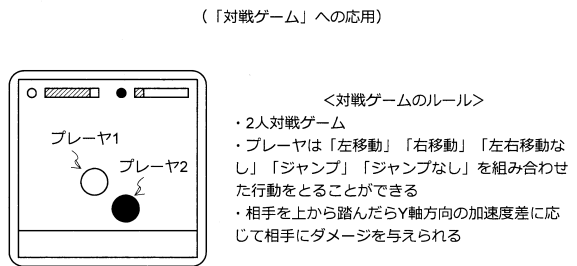
【図53】



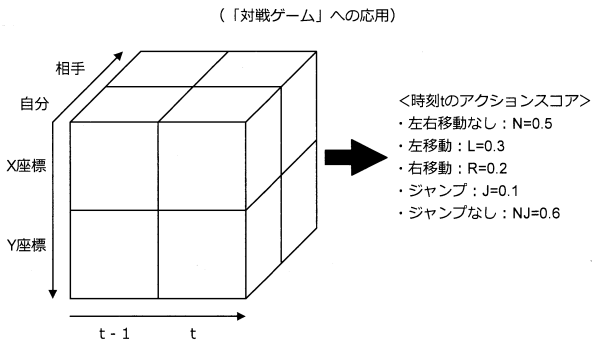
【図54】



【図55】



【図56】



【図57】

(「対戦ゲーム」への応用)

状態Sのバイナリデータ	N	L	R	J	NJ
*as1235asagsdagsd554ads3a	0	1	0	1	0
*lasrigfa543kjadfa674khla	0	0	1	0	1
*687ags87dagsd554ads3a	0	1	0	1	0
*sdagsd55as1235asag4ads3a	0	1	0	0	1
*gsd55as123sda5asag4ads3a	0	0	1	1	0
*s4ad687d55ags87dags3a	0	1	0	0	1
*87dad687d55gs3as4aags	0	1	0	0	1
*5ags8gs3as4ad7da687d5	1	0	0	1	0
⋮	⋮	⋮	⋮	⋮	⋮

時間 ↓

報酬 r = 1 (row 2)

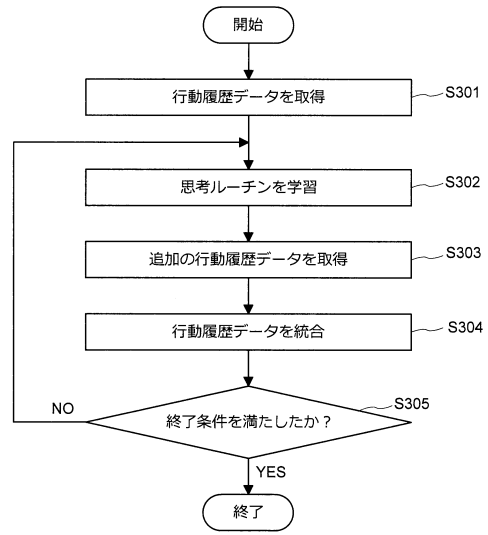
ダメージを与えた (row 4)

報酬 r = 0 (row 7)

ダメージを受けた (row 8)

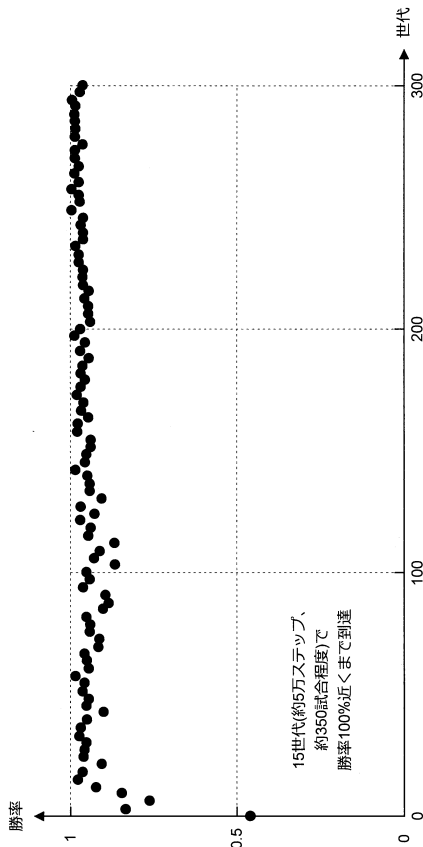
【図58】

(「対戦ゲーム」への応用)



【図59】

(「対戦ゲーム」への応用)



【図60】

(「対戦ゲーム」への応用：複数プレイヤーの同時学習)



【図 6 1】

(「対戦ゲーム」への応用：複数プレイヤーの同時学習)

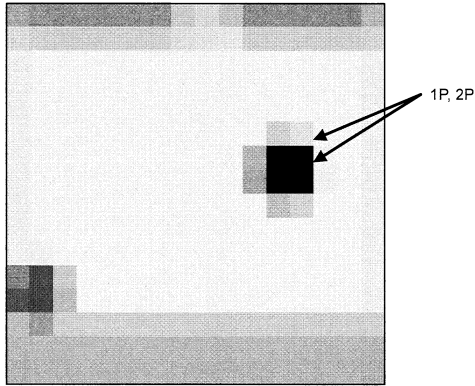
学習後1000試合の評価結果

	対戦勝率	ランダム相手
1P	90.6 %	92.4 %
2P	9.4 %	91.5 %

【図 6 2】

(「対戦ゲーム」への応用)

状態S = ゲーム画面(16x16ブロック)の輝度画像



【図 6 3】

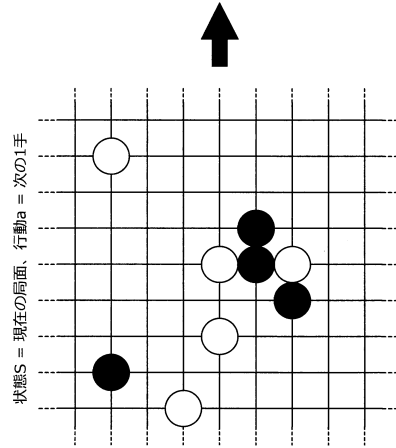
黒="1", 白="-1", 空白="0"

```

... 0000000000 ...
... 0100000-10 ...
... 0000000000 ...
... -1000000000 ...
... 00-10-10000 ...
... 00000110000 ...
... 0001-100000 ...
... 00000000000 ...

```

(「五目並べ」への応用)



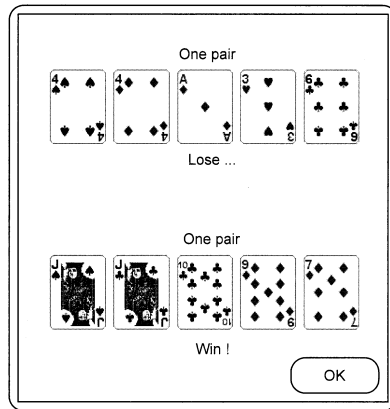
【図 6 4】

(「五目並べ」への応用)



【図 6 5】

(「ポーカー」への応用)



<ポーカーのルール>

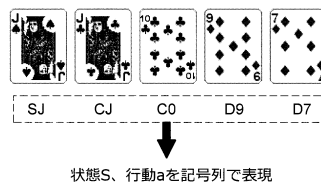
- ・5枚のカードを配る
- ・捨てるカードを選択
- ・役の強い方が勝ち

↓

捨てるカードを決める
思考ルーチンを構築

【図 6 6】

(「ポーカー」への応用)



状態S、行動aを記号列で表現

【 図 6 7 】

(「ポーカー」への応用)

行動履歴データ

状態S：行動aを表す記号列	報酬r
D7SJHKCJS9: SJ HK CJ S9	1
H5DJH0D4DQ: D4	1
S6CQSJSKH0: CQ SK	1
CQH0DAD9D5: HQ D9	1
C2D2H4HQDA: D2 H4	0
⋮	⋮

【 図 6 8 】

(「ロールプレイングゲーム」への応用)

KNIGHT HP: 43 MP: 100 LV: 1	THIEF HP: 64 MP: 100 LV: 1	WITCH HP: 100 MP: 80 LV: 1	PRIEST HP: 100 MP: 90 LV: 1
---	--	--	---

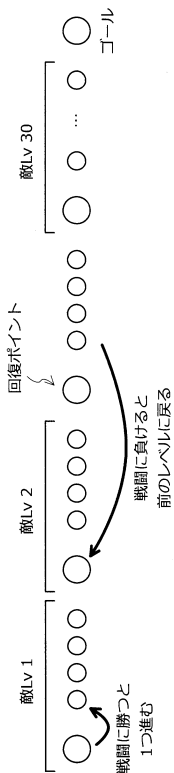
DRAGON1 attacks!
72 damage points for THIEF

<RPGルール>

- ・4人パーティ ・敵3-5匹
- ・状態S = プレーヤに提供される情報 (ex. ステータス、敵の種類など)
- ・行動a = コマンド
(ex. 戦う→ターゲット、魔法→魔法の種類→ターゲット)

【 図 6 9 】

(「ロールプレイングゲーム」への応用：環境について)



- ・戦闘に勝つと生存者で経験値を山分け
- ・経験値が貯まるとレベルアップして強くなる
→ 職業に応じてステータスの値がアップする (ex. 魔法を覚えるなど)
- ・戦闘に5回連続勝つと敵のレベルが1上がり、全回復
- ・全滅すると敵のレベルが1下がり、全回復
- ・敵のレベルが31に達するとクリア

【 図 7 0 】

(「ロールプレイングゲーム」への応用：職業について)

職業	ステータス上昇率 (倍)					魔法を覚えるLv			魔法効果 (倍)		
	最大HP	最大MP	攻撃力	防御力	素早さ	Heal	Fire	Ice	Heal	Fire	Ice
KNIGHT	1.06	1.00	1.06	1.05	1.04				1	1	1
THIEF	1.05	1.00	1.05	1.05	1.06				1	1	1
WITCH	1.04	1.07	1.04	1.04	1.05		0	5	1	1	1
PRIEST	1.05	1.05	1.04	1.05	1.05	0			1	1	1
SLIME	1.05	1.05	1.05	1.05	1.05	10	10	15	1	1	1
UNDEAD	1.05	1.05	1.05	1.06	1.04			15	0	1.5	0.5
DRAGON	1.06	1.05	1.06	1.05	1.05		5	10	1	0.5	1.5
MACHINE	1.05	1.00	1.05	1.05	1.06				1	0	0

【図71】

(「ロールプレイングゲーム」への応用：状態Sに含まれる味方情報)

職業	KNIGHT	THIEF	WITCH	PRIEST	SLIME	UNDEAD	DRAGON	MACHINE
	味方1	味方2	味方3	当てはまるところは"1"、それ以外は"0"				
HP	最大HP	MP	最大MP	攻撃力	防御力	素早さ		

【図72】

(「ロールプレイングゲーム」への応用：状態Sに含まれる敵情報)

職業	KNIGHT	THIEF	WITCH	PRIEST	SLIME	UNDEAD	DRAGON	MACHINE
	敵1	敵2	敵3	当てはまるところは"1"、それ以外は"0"				
累積ダメージ								

【図73】

(「ロールプレイングゲーム」への応用：行動a(生存している味方))

行動者	行動対象	行動の種類			
		攻撃	Heal	Fire	Ice
味方1	対象に"1"、それ以外は"0"				
味方2					
味方3					

【図75】

(「ロールプレイングゲーム」への応用)

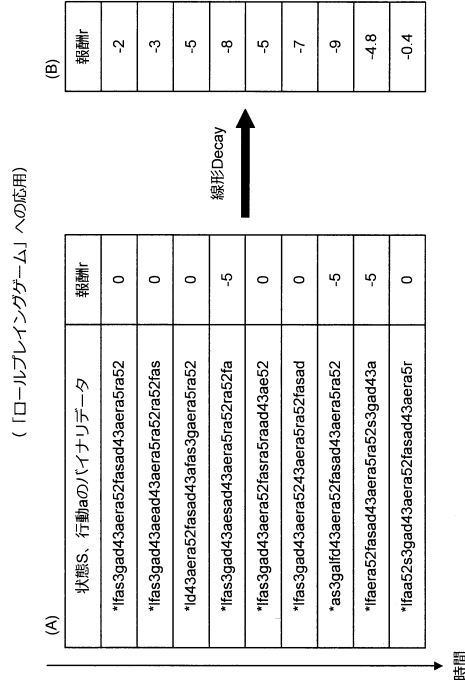
時刻	状態S	行動a	報酬r
0	状態S ₀	敵1を攻撃(a ₀)	0
1	状態S ₁	味方1にHeal(a ₁)	-5 (味方がやられた)
2	状態S ₂	敵2を攻撃(a ₂)	0
3	状態S ₃	敵全体にFire(a ₃)	1 (敵を倒した)

【図74】

(「ロールプレイングゲーム」への応用：行動a(敵))

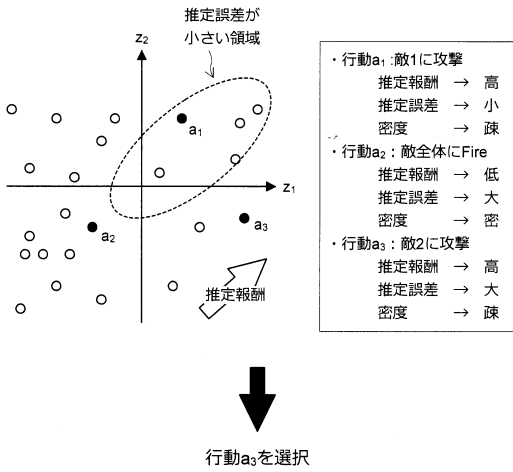
	行動対象
敵1	
敵2	
敵3	

【図 76】

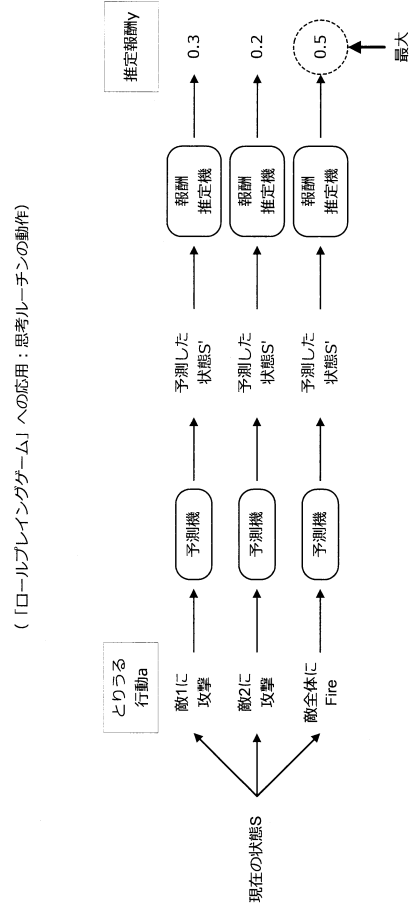


【図 78】

(「ロールプレイングゲーム」への応用：思考ルーチンの動作)

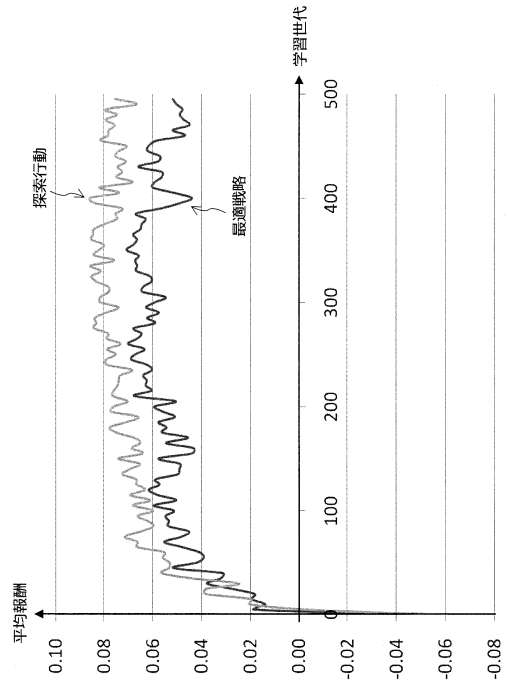


【図 77】

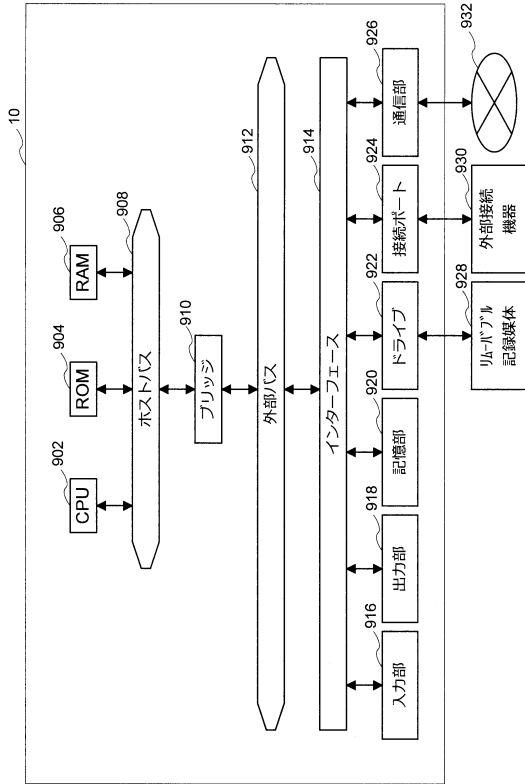


【図 79】

(「ロールプレイングゲーム」への応用：評価結果)



【 図 8 0 】



フロントページの続き

審査官 多賀 実

(56)参考文献 特開2013-058095(JP,A)

木村 元 外1名, "ロボットの強化学習における状態-行動空間の汎化", 日本ロボット学会誌, 社団法人日本ロボット学会, 2004年 3月15日, 第22巻, 第2号, pp. 161-164

堀之内 剛史 外2名, "遺伝的プログラミングを用いたゲームの局面評価関数の学習", 電子情報通信学会技術研究報告, 社団法人電子情報通信学会, 1997年 3月18日, 第96巻, 第594号, pp. 17-24

藤井 叙人 外1名, "戦略型トレーディングカードゲームのための戦略獲得手法", 情報処理学会論文誌, 社団法人情報処理学会, 2009年12月15日, 第50巻, 第12号, pp. 2796-2806

生天目 章 外2名, "エージェント間の相互作用: 望ましい関係性の創発", 計測と制御, 社団法人計測自動制御学会, 2005年12月10日, 第44巻, 第12号, pp. 865-874

Pieter Abbeel 外1名, "Apprenticeship learning via inverse reinforcement learning", ICML '04 Proceedings of the twenty-first international conference on Machine learning, 米国, ACM, 2004年, pp. 1-8

(58)調査した分野(Int.Cl., DB名)

G06N 3/00-99/00

A63F 13/00-13/98