

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

**特許第6145766号  
(P6145766)**

(45) 発行日 平成29年6月14日(2017.6.14)

(24) 登録日 平成29年5月26日(2017.5.26)

(51) Int. Cl.		F I			
<b>G06Q 50/10</b>	<b>(2012.01)</b>	G06Q	50/10		
<b>G06F 19/00</b>	<b>(2011.01)</b>	G06F	19/00	1 2 0	

請求項の数 10 (全 16 頁)

(21) 出願番号	特願2013-66768 (P2013-66768)	(73) 特許権者	503359821 国立研究開発法人理化学研究所 埼玉県和光市広沢2番1号
(22) 出願日	平成25年3月27日(2013.3.27)	(74) 代理人	100120868 弁理士 安彦 元
(65) 公開番号	特開2014-191598 (P2014-191598A)	(72) 発明者	金 成主 埼玉県和光市広沢2番1号 独立行政法人 理化学研究所内
(43) 公開日	平成26年10月6日(2014.10.6)	(72) 発明者	青野 真士 埼玉県和光市広沢2番1号 独立行政法人 理化学研究所内
審査請求日	平成27年11月26日(2015.11.26)	(72) 発明者	行田 悦資 埼玉県和光市広沢2番1号 独立行政法人 理化学研究所内

最終頁に続く

(54) 【発明の名称】 解探索システム、解探索プログラム

(57) 【特許請求の範囲】

## 【請求項1】

確率分布に基づいて結果を出力する2以上の被検対象のうち最良の結果の出力が期待される被検対象を探索する解探索システムにおいて、

上記出力された結果の蓄積に基づく今までの戦績を上記被検対象毎にそれぞれ求め、上記各被検対象の戦績を被検対象全体の戦績との関係においてその優劣を比較する戦績優劣比較手段と、

上記戦績優劣比較手段により比較された戦績の優劣と、上記被検対象から出力された直近の結果とに基づいて、計量変数を増加又は減少させるように制御することを当該被検対象毎に行う制御手段と、

上記計量変数が閾値を超えた上記被検対象に対して結果の出力を指示する出力指示手段とを備え、

上記出力指示手段は、上記結果の出力の指示の繰り返しを経て最終的に最も上記結果の出力の指示が行われている1以上の被検対象を、探索解として特定すること

を特徴とする解探索システム。

## 【請求項2】

それぞれ設定された確率分布に基づいて結果を出力する3以上の被検対象のうち最良の結果の出力が期待される被検対象の組み合わせを探索し、

上記出力指示手段は、上記結果の出力の指示の繰り返しを経て最終的に最も上記結果の出力の指示が行われている被検対象の組み合わせを、探索解として特定すること

を特徴とする請求項 1 記載の解探索システム。

【請求項 3】

上記出力指示手段による出力指示に応じて、時系列的に確率分布が変化する被検対象の組み合わせを探索する上で、2 以上の当該解探索システムを用いて解探索を行うことを特徴とする請求項 1 記載の解探索システム。

【請求項 4】

上記戦績優劣比較手段は、一の被検対象がより優れた結果を出力した場合及び他の被検対象がより劣る結果を出力した場合に、上記一の被検対象における今までの戦績をより優れた側に向きさせ、一の被検対象がより劣った結果を出力した場合及び他の被検対象がより優れた結果を出力した場合に、上記一の被検対象における今までの戦績をより劣る側に下降させること

を特徴とする請求項 1 ~ 3 のうち何れか 1 項記載の解探索システム。

【請求項 5】

上記戦績優劣比較手段は、上記各被検対象の戦績と被検対象全体の戦績の平均との差分を内部リソース値とし、

上記制御手段は、

上記内部リソース値が正で、上記被検対象から出力された直近の結果がより優れたものである場合には、上記計量変数を増加させ、

上記内部リソース値が正で、上記被検対象から出力された直近の結果がより劣るものである場合には、上記計量変数をそのままにし、

上記内部リソース値が 0 で、上記被検対象から出力された直近の結果がより優れたものである場合には、上記計量変数を増加させ、

上記内部リソース値が 0 で、上記被検対象から出力された直近の結果がより劣るものである場合には、上記計量変数を下降させ、

上記内部リソース値が負で、上記被検対象から出力された直近の結果がより優れたものである場合には、上記計量変数をそのままにし、

上記内部リソース値が負で、上記被検対象から出力された直近の結果がより劣るものである場合には、上記計量変数を下降させること

を特徴とする請求項 1 ~ 4 のうち何れか 1 項記載の解探索システム。

【請求項 6】

上記制御手段は、上記被検対象毎に割り当てられる計量変数の合計が一定となるように制御すること

を特徴とする請求項 1 ~ 5 のうち何れか 1 項記載の解探索システム。

【請求項 7】

確率分布に基づいて結果を出力する 2 以上の被検対象のうち最良の結果の出力が期待される被検対象を探索する解探索プログラムにおいて、

上記出力された結果の蓄積に基づく今までの戦績を上記被検対象毎にそれぞれ求め、上記各被検対象の戦績を被検対象全体の戦績との関係においてその優劣を比較する戦績優劣比較ステップと、

上記戦績優劣比較ステップにより比較された戦績の優劣と、上記被検対象から出力された直近の結果とに基づいて、計量変数を増加又は減少させるように制御することを当該被検対象毎に行う制御ステップと、

上記計量変数が閾値を超えた上記被検対象に対して結果の出力を指示する出力指示ステップとを有し、

上記出力指示ステップでは、上記結果の出力の指示の繰り返しを経て最終的に最も上記結果の出力の指示が行われている 1 以上被検対象を、探索解として特定すること

をコンピュータに実行させることを特徴とする解探索プログラム。

【請求項 8】

上記戦績優劣比較ステップでは、一の被検対象がより優れた結果を出力した場合及び他の被検対象がより劣る結果を出力した場合に、上記一の被検対象における今までの戦績を

より優れる側に向上させ、一の被検対象がより劣った結果を出力した場合及び他の被検対象がより優れた結果を出力した場合に、上記一の被検対象における今までの戦績をより劣る側に下降させること

を特徴とする請求項 7 記載の解探索プログラム。

【請求項 9】

上記戦績優劣比較ステップでは、上記各被検対象の戦績と被検対象全体の戦績の平均との差分を内部リソース値とし、

上記制御ステップでは、

上記内部リソース値が正で、上記被検対象から出力された直近の結果がより優れたものである場合には、上記計量変数を増加させ、

上記内部リソース値が正で、上記被検対象から出力された直近の結果がより劣るものである場合には、上記計量変数をそのままにし、

上記内部リソース値が 0 で、上記被検対象から出力された直近の結果がより優れたものである場合には、上記計量変数を増加させ、

上記内部リソース値が 0 で、上記被検対象から出力された直近の結果がより劣るものである場合には、上記計量変数を下降させ、

上記内部リソース値が負で、上記被検対象から出力された直近の結果がより優れたものである場合には、上記計量変数をそのままにし、

上記内部リソース値が負で、上記被検対象から出力された直近の結果がより劣るものである場合には、上記計量変数を下降させること

を特徴とする請求項 7 又は 8 記載の解探索プログラム。

【請求項 10】

上記制御ステップでは、上記被検対象毎に割り当てられる計量変数の合計が一定となるように制御すること

を特徴とする請求項 7 ~ 9 のうち何れか 1 項記載の解探索プログラム。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、組み合わせ報酬最大化問題の解を高速かつ効率的に導く上で好適な解探索システム、解探索プログラムに関するものである。

【背景技術】

【0002】

従来より、期待値を最大化する解を探索する問題の代表例として、バンディット問題がある。このバンディット問題とは、貰える合計報酬の期待値を最大化することを目的とし、プレイヤーは  $n$  種類の異なる行動選択肢から一つの選択肢を選択する動作を繰り返す。各選択の後には毎回、選択した行動に依存する確率分布から選ばれた結果がプレイヤーの報酬として与えられる。

【0003】

仮に複数のスロットマシンが存在し、各スロットマシンのレバーを引くことにより、ある確率分布の下でコイン（報酬）がもらえるものとする。このコインが出る確率分布（当選確率）がスロットマシン毎に異なる場合であって、かつプレイヤーはその当選確率が分からない場合を考えてみる。このとき、各スロットマシンの当選確率を知る最も一般的な方法としては、とりあえず各スロットマシンを多数回に亘り順にプレイし、実際に最も報酬が大きかったスロットマシンが、最も当選確率が高いものと判断する。

【0004】

しかしながら、かかる方法では、実際に最も当選確率の高いスロットマシンを特定する上で相当の回数に亘りスロットマシンをプレイしなければならず、結果として多くの投資が必要となる。このため、各スロットマシンにおける当選確率を調べる上で極力投資を少なくしつつ、効率的に解を探索できるアルゴリズムを考える必要が出てくることが分かる。

10

20

30

40

50

## 【 0 0 0 5 】

かかる場合には、上述のような貰える合計報酬の期待値を最大化するバンディット問題に置き換えて、これを解くことができる（例えば、非特許文献1参照。）。特にこのバンディット問題の中で、n種類の異なる行動選択肢から最良の結果の出力が期待される組み合わせを選択する、いわゆる組み合わせバンディット問題も近年において注目されている。この組み合わせバンディット問題では、複数台のスロットマシンの中からより高配当が期待できるスロットマシンの組み合わせを選択する場合のみならず、例えばコグニティブ無線通信においてデータ伝送量を最大化できるチャネルの最適な組み合わせの選択、インターネット広告においてクリック数を最大化できる広告の最適な組み合わせ、更には最も投資リターンの大きい金融商品のポートフォリオの選択等、様々な分野においてニーズがある。このような応用例の場合には、より一般的な組み合わせ報酬最大化問題になる。つまり、プレイヤーが多数で、それぞれのプレイヤーの選択に依存して（例えばペイオフ行列によって）各プレイヤーの報酬量が決定される。しかし、本明細書では簡単のために、各スロットマシンが独立な場合の、組み合わせ報酬最大化問題（特に2つの組み合わせ）について例を挙げて説明する。

10

## 【 先行技術文献 】

## 【 非特許文献 】

## 【 0 0 0 6 】

【 非特許文献 1 】 S.J.Kim,M.Aono,M.Hara, BioSystems 101,29 36 (2010)

## 【 発明の概要 】

20

## 【 発明が解決しようとする課題 】

## 【 0 0 0 7 】

しかしながら、従来において、組み合わせバンディット問題の解を自動的に探索して求めるためのアルゴリズムが特段提案されていなかった。情報量が増大の一途を辿る昨今において、大量の情報から高速かつ効率的に、組み合わせバンディット問題の解を求めるための社会的要請が高くなると考えられるが、これについて特段の解決策が提案されていないのが現状であった。

## 【 0 0 0 8 】

本発明は、上述した問題点に鑑みて案出されたものであり、その目的とするところは、高速かつ効率的に組み合わせバンディット問題の解を求めることが可能な解探索システム、解探索プログラムを提供することにある。

30

## 【 課題を解決するための手段 】

## 【 0 0 0 9 】

本発明を適用した解探索システムは、上述した課題を解決するために、確率分布に基づいて結果を出力する2以上の被検対象のうち最良の結果の出力が期待される被検対象を探索する解探索システムにおいて、上記出力された結果の蓄積に基づく今までの戦績を上記被検対象毎にそれぞれ求め、上記各被検対象の戦績を被検対象全体の戦績との関係においてその優劣を比較する戦績優劣比較手段と、上記戦績優劣比較手段により比較された戦績の優劣と、上記被検対象から出力された直近の結果とに基づいて、計量変数を増加又は減少させるように制御することを当該被検対象毎に行う制御手段と、上記計量変数が閾値を超えた上記被検対象に対して結果の出力を指示する出力指示手段とを備え、上記出力指示手段は、上記結果の出力の指示の繰り返しを経て最終的に最も上記結果の出力の指示が行われている1以上の被検対象を、探索解として特定する解探索システムことを特徴とする。

40

## 【 0 0 1 0 】

本発明を適用した組み合わせ探索プログラムは、確率分布に基づいて結果を出力する2以上の被検対象のうち最良の結果の出力が期待される被検対象を探索する解探索プログラムにおいて、上記出力された結果の蓄積に基づく今までの戦績を上記被検対象毎にそれぞれ求め、上記各被検対象の戦績を被検対象全体の戦績との関係においてその優劣を比較する戦績優劣比較ステップと、上記戦績優劣比較ステップにより比較された戦績の優劣と、

50

上記被検対象から出力された直近の結果とに基づいて、計量変数を増加又は減少させるように制御することを当該被検対象毎に行う制御ステップと、上記計量変数が閾値を超えた上記被検対象に対して結果の出力を指示する出力指示ステップとを有し、上記出力指示ステップでは、上記結果の出力の指示の繰り返しを経て最終的に最も上記結果の出力の指示が行われている1以上被検対象を、探索解として特定することをコンピュータに実行させることを特徴とする。

【発明の効果】

【0012】

上述した構成からなる本発明によれば、被検対象から出力された今までの戦績の優劣と、被検対象から出力された直近の結果とに基づいて、計量変数を増加又は減少させ、この計量変数の値に応じて結果の出力を指示するか否かを決定する。そして、結果の出力の指示の繰り返しを経て最終的に最も結果の出力の指示が行われている被検対象の組み合わせを、探索すべき組み合わせとして特定する。これにより、組み合わせバンディット問題の解を自動的に探索して求めることが可能となり、情報量が増大の一途を辿る昨今において、大量の情報から高速かつ効率的に、組み合わせバンディット問題の解を求めることが可能となる。

【図面の簡単な説明】

【0013】

【図1】本発明を適用した解探索システムの全体構成を示す図である。

【図2】計量変数について説明するための図である。

【図3】本発明を適用した解探索システムの他の全体構成例を示す図である。

【図4】本発明を適用した解探索システムを、コンピュータプログラムで実現した場合における実施例を示す図である。

【図5】被検対象の当選確率が、0.2、0.5、0.8の3つのサンプルについての組み合わせを選択するシミュレーションを行う例を示す図である。

【発明を実施するための形態】

【0014】

以下、本発明を適用した解探索システムについて図面を参照しながら詳細に説明をする。

【0015】

図1は、本発明を適用した解探索システム解探索システム1の全体構成を示している。この解探索システム1は、2以上の被検対象5a~5dのうち最良の結果の出力が期待される被検対象5の組み合わせを探索するシステムである。解探索システム1は、戦績優劣比較部2と、この戦績優劣比較部2に接続された制御部3a~3dと、制御部3a~3dに接続された出力指示部4a~4dとを備えている。

【0016】

ちなみに、この制御部3a~3d、出力指示部4a~4dは、それぞれ被検対象5の数と同等となるように設けられるものであり、図1の例では4個の被検対象5からなるため、これら制御部3、出力指示部4も4個ずつで構成される。

【0017】

被検対象5は、それぞれ設定された確率分布に基づいて結果を出力する対象物である。例えば、スロットマシンやパチンコの台のように、設定された確率分布に基づいてコインという結果物を出力するものであってもよい。また、コグニティブ無線通信は、各チャンネルのデータ伝送量の大小は、その都度変化するものであるが、これについてもある時点において設定された確率分布で表現することができる。このようなコグニティブ無線通信において任意のチャンネルを選択した場合に、実際の"データ伝送量"という結果物を出力する。また、インターネット広告についても、掲載すべき広告のクリック数の大小は確率分布で表すことが可能となり、"実際のクリック数"という結果物はその確率分布に基づいて算出することが可能となる。また、金融商品については、その将来的な投資リターンも確率分布で表すことができ、"実際の投資リターン"という結果物も当該確率分布に基

づいて表される。

【0018】

このように、被検対象5は、出力する結果を確率分布に変換することが可能なあらゆる事象、物、システム、プログラムやアルゴリズムを含む概念である。ちなみに、この被検対象5において出力される結果の確率分布は、通常の正規分布、ガウシアン分布のみならず、離散的な分布であってもよいし、2項分布で構成されていてもよい。ちなみに、この被検対象5の確率分布は、この解探索システム1のユーザにとって未知のものとなっている。ユーザは、これらの被検対象のうち、最良の結果の出力が期待される被検対象の組み合わせを探索するためにこの解探索システム1を使用することとなる。

【0019】

被検対象5 a ~ 5 dは、それぞれ設定された確率分布に基づいて結果1を出力する。このとき、被検対象5 aから出力される結果を結果1とし、被検対象5 bから出力される結果を結果1とし、被検対象5 cから出力される結果を結果1とし、被検対象5 dから出力される結果を結果1とする。出力された結果1は、戦績優劣比較部2へと送信されると共に、制御部3 aへ送信される。出力された結果1は、戦績優劣比較部2へと送信されると共に、制御部3 bへ送信される。出力された結果1は、戦績優劣比較部2へと送信されると共に、制御部3 cへ送信される。出力された結果1は、戦績優劣比較部2へと送信されると共に、制御部3 dへ送信される。

【0020】

戦績優劣比較部2は、被検対象5 a ~ 5 dから出力される結果1 ~ 1を受信し、これを記憶する。この戦績優劣比較部2は、被検対象5 a ~ 5 dから結果1 ~ 1を受信する都度、順次記憶しておくことで、結果を蓄積する。そして、この戦績優劣比較部2は、被検対象5 a ~ 5 d毎に、出力された結果の蓄積に基づく今までの戦績をそれぞれ求める。ここでいう戦績とは、被検対象5から出力される結果がより優れているのか、或いはより劣っているのかを示すあらゆるデータを示すものである。被検対象5がスロットマシンであれば、コインがどの程度出たかを示すものであってもよいし、被検対象5がインターネット広告であれば、クリックがどの程度行われたかを示すデータであってもよい。また、この戦績優劣比較部2は、各被検対象5の戦績を被検対象5全体の戦績との関係においてその優劣を比較する。ここで被検対象5 aについての被検対象5全体の優劣を $s_1$ とし、被検対象5 bについての被検対象5全体の優劣を $s_2$ とし、被検対象5 cについての被検対象5全体の優劣を $s_3$ とし、被検対象5 dについての被検対象5全体の優劣を $s_4$ とする。これら優劣 $s_1$ は、制御部3 aに、優劣 $s_2$ は、制御部3 bに、優劣 $s_3$ は、制御部3 cに、優劣 $s_4$ は、制御部3 dにそれぞれ送られる。

【0021】

制御部3 a ~ 3 dは、それぞれ戦績優劣比較部2から、戦績の優劣に関する情報 $s_1$  ~  $s_4$ がそれぞれ入力されるとともに、被検対象5 a ~ 5 dから直近の結果1 ~ 1がそれぞれ入力される。制御部3 a ~ 3 dは、それぞれ入力された戦績の優劣に関する情報 $s_1$  ~  $s_4$ と、被検対象から出力された直近の結果とに基づいて計量変数を増加又は減少させるように制御する。

【0022】

図2は、この計量変数のイメージを説明するための図である。計量変数 $x$ は、各被検対象5に対してそれぞれ個別に割り当てられたパラメータであり、被検対象5を選択するか否かの判断する上での基準となるものである。この計量変数 $x$ が高いほど、その被検対象5を選択する可能性が高くなる。一方、この計量変数 $x$ が低いほどその被検対象5を選択しない可能性が高くなる。

【0023】

このような計量変数 $x$ の増減を制御するのが、制御部3である。つまり制御部3 aは、結果1と優劣 $s_1$ とに基づいて計量変数 $x$ の増減を制御し、制御部3 bは、結果1と優劣 $s_2$ とに基づいて計量変数 $x$ の増減を制御し、制御部3 cは、結果1と優劣 $s_3$ とに基づいて計量変数 $x$ の増減を制御し、制御部3 dは、結果1と優劣 $s_4$ とに基づ

10

20

30

40

50

いて計量変数  $x_i^t$  の増減を制御する。これら制御部 3 による計量変数  $x_i^t$  の制御は、あくまで結果 1 と、優劣  $s^t$  に基づくものであればいかなるものであってもよい。

【0024】

出力指示部 4 a ~ 4 d は、計量変数を監視し、予め設定した閾値を超えるか否かを判別する。そして、計量変数が閾値を超えた旨を判別した場合には、被検対象 5 に対して結果の出力を指示する。この図 2 の例では、閾値を超えた計量変数は、 $x_2^t$  と  $x_4^t$  である。このため、 $x_2^t$  を監視する出力指示部 4 b と、 $x_4^t$  を監視する出力指示部 4 d は、被検対象 5 b、5 d に対して結果の出力を指示する。

【0025】

このようにして、被検対象 5 を中心として、戦績優劣比較部 2、制御部 3、出力指示部 4 の順でいわゆるフィードバック制御が行われる。このような戦績優劣比較部 2、制御部 3、出力指示部 4 からなる解探索システム 1 は、例えば、アナログ回路、デジタル回路を始めとしたいかなるデバイスで具現化されるものであってもよい。ちなみに、回路として具現化される場合には、FPGA(field programmable gate array)に基づいて、構成を設定するようにしてもよい。また、本発明はプログラムで具現化されるものであってもよい。かかる場合には、戦績優劣比較部 2 は、これと同様の処理を実行する戦績優劣比較ステップに、制御部 3 は、これと同様の処理を実行する制御ステップに、出力指示部 4 は、これと同様の処理を実行する出力指示ステップとして具現化されることとなる。これに加えて、このようなプログラムに基づいて動作するハードウェア（例えば、パーソナルコンピュータ、各種携帯情報端末等）を介して具現化されるものであってもよい。

【0026】

次に、本発明を適用した解探索システム 1 による組み合わせ探索動作について説明をする。

【0027】

まず、被検対象 5 a ~ 5 d のいくつかについて結果の出力が行われる。この結果の出力が指示されるのは上述した出力指示部 4 により出力指示が行われたものに限る。

【0028】

被検対象 5 a ~ 5 d から出力された結果 1 は、それぞれ戦績優劣比較部 2 へ送られると共に、制御部 3 へと送られる。戦績優劣比較部 2 では、この送られてきた結果 1 に基づいて具体的に以下の処理動作を行う。

【0029】

以下の式 (1) における  $q_i^t$  は、それぞれの被検対象 5 a ~ 5 d における戦績を示す指数である。この戦績指数  $q_i^t$  の  $t$  は、この解探索システム 1 におけるフィードバック回数を示している。また、 $i$  は、図 1 に示す下付きの数字に対応するものであり、それぞれの被検対象 5 に対応するものである。つまり  $i$  毎に、換言すれば被検対象 5 毎に、この戦績指数  $q_i^t$  を求めていくこととなる。

【0030】

【数 1】

$$q_i^t = \alpha \cdot q_i^{t-1} + \sum_{j \in \Lambda(i)} \mu \cdot (\rho_i^t \cdot \rho_j^t + \omega \cdot \sum_{\langle k, k' \rangle \in I_2 \setminus \{\langle i, j \rangle\}} \pi_k^t \cdot \pi_{k'}^t)$$

..... (1)

【0031】

戦績指数  $q_i^t$  は、フィードバックの回数  $t$  が増加するにつれて順次更新される。そして、この戦績指数  $q_i^t$  は、その前のフィードバック回数  $t - 1$  において求めた戦績指数  $q_i^{t-1}$  に加えて、以降の項を足すことで表示される。また、この (1) 式においては忘却パラメータであり、必要に応じて設定される。

【0032】

この以降の項において、 $\mu$  は係数である。また、 $\rho_i^t$  は、より優れた結果を出力

10

20

30

50

した場合を意味している。ここで  $\rho$  は、これから戦績指数  $q_i$  を求めようとする被検対象 5 が、より優れた結果を出力した場合に " 1 " となり、より劣る結果を出力した場合には、" 0 " となる。 $\pi$  は、 $j$  に相当する他の被検対象 5 が、より優れた結果を出力した場合には " 1 " となり、より劣る結果を出力した場合には、" 0 " となる。本明細では、例として、2つの組み合わせを選択するためのシステムであることから  $\rho$  は2変数となっている。

【0033】

表1は、 $\rho$ 、 $\pi$  の取り得る値を示している。

【0034】

【表1】

	優	劣	非出力
$\rho$	1	0	0
$\pi$	0	1	0

【0035】

式(1)の括弧内は、 $\rho \cdot \pi$  のように乗算で表されることから、 $\rho$  と  $\pi$  の双方が1である場合のみプラスになる。

【0036】

また、 $\pi$  は、 $i$ 、 $j$  以外の他の被検対象 5 が、より劣った結果を出力した場合に " 1 " となり、より優れた結果を出力した場合に " 0 " となる。

【0037】

即ち、この(1)式の括弧内は、その  $q_i$  を求める被検対象 5 が、より優れた結果を出力した場合と、当該  $q_i$  を求める被検対象 5 以外の他の被検対象 5 が、より劣った結果を出力した場合に、その数値が上昇することになっている。仮に被検対象 5 がスロットマシンである場合には、当該  $q_i$  を求めるスロットマシンが当選した場合には、(1)式の括弧内の数値を上昇させ、当該  $q_i$  を求めるスロットマシンが落選した場合には、(1)式の括弧内の数値を変化させない。また、当該  $q_i$  を求めるスロットマシン以外が当選した場合には、(1)式の括弧内の数値を変化させず、当該  $q_i$  を求めるスロットマシン以外が落選した場合には、(1)式の括弧内の数値を上昇させる。そして、(1)式の括弧内の数値が上昇するにつれて戦績指数  $q_i$  が上昇し、より優れた戦績になる。

【0038】

このように、本発明では、 $q_i$  を求める一の被検対象 5 が、より優れた結果を出力した場合及び他の被検対象がより劣る結果を出力した場合に、当該一の被検対象 5 における今までの戦績  $q_i$  をより優れる側に向上させる。また、 $q_i$  を求める一の被検対象 5 が、より劣った結果を出力した場合及び他の被検対象がより優れた結果を出力した場合に、当該一の被検対象 5 における今までの戦績  $q_i$  を変化させない。戦績優劣比較部 2 は、このような制御を行うものであってもよく、上述した(1)式に基づく制御を行う場合に限定されるものではない。

【0039】

また、何れの出力結果が優れており、何れの出力結果が劣っているかについては、いかなる基準の下で判断するようにしてもよい。上述したスロットマシンの例では、コインが出るか否かで優劣を決める場合に限定されるものではなく、コインの枚数や種別に応じて優劣を決めるようにしてもよい。また、この優劣についても、優れているか、或いは劣っているかの2段階で設定される場合に限定されるものではなく、3段階以上で優劣を評価するようにしてもよい。ちなみに、3段階以上で優劣をランク分けする場合においても、一の被検対象 5 がより上位ランクであるほど  $q_i$  を上昇させ、他の被検対象 5 がより下位ランクであるほど  $q_i$  を下降させるように調整を行う。

【0040】

10

20

30

40

50



戦績優劣比較部 2 は、このようにして、各被検対象 5 a ~ 5 d についてそれぞれ戦績指数  $q_i^t \sim q_i^t$  を求めた後、その戦績指数  $q_i^t \sim q_i^t$  を被検対象全体 5 の戦績との関係においてその優劣を比較する。そして、この優劣の比較結果としての内部リソース値  $s_i^t$  をそれぞれ出力する。

【 0 0 4 1 】

かかる場合に、( 2 ) 式に基づいてその優劣を比較するようにしてもよい。

【 0 0 4 2 】

【 数 2 】

$$s_i^t = x_0^t + q_i^{t-1} - \text{Mean}_{k \in \Lambda(i)} \{q_k^{t-1}\}$$

..... ( 2 )

【 0 0 4 3 】

式 ( 2 ) は、一の被検対象 5 の戦績  $q_i^t$  と被検対象 5 全体の戦績  $q_i^t$  の平均との差分を、内部リソース値  $s_i^t$  としている。例えば、被検対象 5 a についての内部リソース値  $s_i^t$  を求める場合には、その戦績  $q_i^t$  と、全被検対象の戦績  $q_i^t$ 、 $q_i^t$ 、 $q_i^t$  の平均との差分を求める。ちなみに、( 2 ) 式の右項において  $x_i^t$  はあくまで調整値であり、必須のものではない。この ( 2 ) 式に基づいて、内部リソース値  $s_i^t$  を判断する場合に限定されるものではない。内部リソース値  $s_i^t$  は、それに対応する一の被検対象 5 の戦績  $q_i^t$  が、他の被検対象 5 との間で相対的に優れているか否かを示すものであれば、いかなる計算式で、或いはいかなる方法で、これを評価するようにしてもよい。戦績優劣比較部 2 は、これら計算した  $s_i^t \sim s_i^t$  をそれぞれ、制御部 3 a ~ 3 d へ出力する。

20

【 0 0 4 4 】

制御部 3 a ~ 3 d は、これら  $s_i^t \sim s_i^t$  並びに被検対象 5 a ~ 5 d から直近の結果  $l_i^t \sim l_i^t$  がそれぞれ入力された場合に、例えば、以下の制御を行う。

【 0 0 4 5 】

表 2 は、この制御部 3 による制御を行う上で参酌するテーブルの例を示している。

【 0 0 4 6 】

【 表 2 】

30

	$s_i^t > 0$	$s_i^t = 0$	$s_i^t < 0$
$l_i^t = -1$	1	1	0
$l_i^t = 1$	0	-1	-1

【 0 0 4 7 】

この表 2 によれば、制御部 3 は、内部リソース値  $s_i^t$  と、被検対象 5 から出力された直近の結果  $l_i^t$  とより形成されるマトリクスに基づいて、制御方法を決定する。ここで  $l_i^t = -1$  は、より優れた結果が出力された場合を意味している。また、 $l_i^t = 1$  は、より劣った結果が出力された場合を意味している。ちなみに、この  $l_i^t$  の優劣の基準についてもいかなるものであってもよい。

40

【 0 0 4 8 】

また、表中の数値 ( - 1、0、1 ) の意味については、先ず " - 1 " は、計量変数  $x$  を成長速度を減少させるように制御する。また、" 0 " は、計量変数  $x$  の成長速度を特に増減させないことを意味し、" 1 " は、計量変数  $x$  の成長速度を増加させるように制御する。

【 0 0 4 9 】

ちなみに以下の例では、 $x_i^t$  は、 $x_i$  の値によって決まるものと仮定している。

【 0 0 5 0 】

50

制御部 3 は、この表 2 に基づいて具体的に以下の制御を行う。

【 0 0 5 1 】

内部リソース値  $s_i$  が正で、被検対象 5 から出力された直近の結果  $1_i$  がより優れたものである場合 ( $= -1$ ) には、" 1 " であることから計量変数  $x_i$  の成長速度を増加させる。内部リソース値  $s_i$  が正で、被検対象 5 から出力された直近の結果  $1_i$  がより劣るものである場合 ( $= 1$ ) には、" 0 " であることから計量変数  $x_i$  の成長速度を増加させることなくそのままにする。内部リソース値  $s_i$  が 0 で、被検対象 5 から出力された直近の結果  $1_i$  がより優れたものである場合 ( $= -1$ ) には、" 1 " であることから計量変数計量変数  $x_i$  の成長速度を増加させる。内部リソース値  $s_i$  が 0 で、被検対象 5 から出力された直近の結果  $1_i$  がより劣るものである場合 ( $= 1$ ) には、" -1 " であることから計量変数計量変数  $x_i$  の成長速度を減少させる。内部リソース値  $s_i$  が負で、被検対象 5 から出力された直近の結果  $1_i$  がより優れたものである場合 ( $= -1$ ) には、" 0 " であることから計量変数計量変数  $x_i$  の成長速度をそのままにし、内部リソース値  $s_i$  が負で、被検対象 5 から出力された直近の結果  $1_i$  がより劣るものである場合 ( $= 1$ ) には、" -1 " であることから計量変数計量変数  $x_i$  の成長速度を減少させる。

10

【 0 0 5 2 】

このようにして制御部 3 は、戦績優劣比較部 2 により比較された戦績の優劣に基づく内部リソース値  $s_i$  と、被検対象 5 から出力された直近の結果  $1_i$  とに基づいて、計量変数  $x_i$  を増加又は減少させるように制御する。ちなみに表 2 中の数値はあくまで一例であり、内部リソース値  $s_i$  と直近の結果  $1_i$  とに基づくものであればいかなる数値であってもよい。

20

【 0 0 5 3 】

また、上述した例では、内部リソース値  $s_i$  を 3 段階で表示し、直近の結果  $1_i$  を 2 段階で表示し、合計  $2 \times 3$  のマトリックスで表示しているが、これに限定されるものではなく、 $s_i$ 、 $1_i$  とともに最低 2 段階であれば何段階でランク分けするようにしてもよい。

【 0 0 5 4 】

さらに、この制御部 3 は、このようなマトリックス表で制御方法を規定する場合に限定されるものではなく、内部リソース値  $s_i$  と直近の結果  $1_i$  とに基づくものであれば他のいかなる方法に基づいて計量変数  $x_i$  の増減を制御するようにしてもよい。具体的には、内部リソース値  $s_i$  と直近の結果  $1_i$  とを変数とした所定の演算式に従って制御方法を定めるようにしてもよい。

30

【 0 0 5 5 】

このようにして制御部により計量変数  $x_i$  の増減を制御する結果、図 2 に示すように、ある被検対象 5 に対する計量変数  $x_i$  は閾値未満であり、ある被検対象 5 に対する計量変数  $x_i$  は閾値以上となる。出力指示部 4 は、これら軽量パラメータ  $x_i$  をそれぞれ閾値との関係で比較し、閾値を超えた計量変数  $x_i$  に応じた被検対象 5 に対してのみ、結果の出力を指示する。結果の出力を指示された被検対象 5 は、新たな結果をそれぞれの確率分布に基づいて出力することとなる。

【 0 0 5 6 】

本発明を適用した解探索システム 1 では、これらの動作を繰り返し実行する。その結果、出力指示部 4 により出力指示が出される被検対象 5 の組み合わせが徐々に一定化されてくる。最終的には、この出力指示が出される被検対象 5 は、一定のものに収束される。この収束されてくる被検対象 5 の組み合わせが、解探索システム 1 により選択されてくる解となる。

40

【 0 0 5 7 】

ちなみに、解探索システム 1 により解を求める上で、被検対象 5 a ~ 5 d 毎に割り当てられる計量変数  $x_1 \sim x_4$  と、 $x_i$  との合計が一定となるように制御することで、より探索精度を高めることが可能となる。即ち、 $x_1 + x_2 + x_3 + x_4 = \text{一定}$ 、としておくことにより、計量変数  $x_i$  が全体的に大きくなってしまい、閾値との間でバランスが取れなくなるのを防止することができ、これが探索精度の向上につながる事となる。

50

## 【 0 0 5 8 】

$x_i$ は、他の計量変数  $x$  の値により影響を受ける変数である。上述した処理動作を繰り返し進めていく結果、上述した(2)式において、当初は  $x_i$  の項が支配的になり、 $q_i$  以後の項(一の被検対象5の戦績  $q_i$  と被検対象5全体の戦績  $q_i$  の平均との差分)については、あまりこの  $s_i$  を決定する上で大きな影響を与えるものではない。その結果、この  $s_i$  は、戦績  $q_i$  と戦績  $q_i$  の平均との差分に影響を受けることなく、自由度が高くなり、その分ランダムな値を取りやすくなる。その結果、ランダムな値を取りやすくなる  $s_i$  を介して様々な出力指示 4 a ~ 4 d が行われ、様々な解を探索することが可能となる。

## 【 0 0 5 9 】

これに対して、 $t$  が大きくなるにつれ、換言すればフィードバック回数が多くなるにつれ、徐々に  $q_i$  以後の項(一の被検対象5の戦績  $q_i$  と被検対象5全体の戦績  $q_i$  の平均との差分)が大きくなる。その結果、この  $s_i$  を決定する上で、 $s_i$  よりも、戦績  $q_i$  と戦績  $q_i$  の平均との差分値が支配的になってくる。そして、次回に出力指示 4 a ~ 4 d が行われるものについては、戦績  $q_i$  と戦績  $q_i$  の平均との差分値がより影響を受けるものとなる。

## 【 0 0 6 0 】

このようにして最終的に出力指示が行われるものは、戦績  $q_i$  と戦績  $q_i$  の平均との差分値による影響を受けるものに収束されてくる。

## 【 0 0 6 1 】

上述した構成からなる本発明によれば、被検対象5から出力された今までの戦績の優劣と、被検対象5から出力された直近の結果とに基づいて、計量変数を増加又は減少させ、この計量変数の値に応じて結果の出力を指示するか否かを決定する。そして、結果の出力の指示の繰り返しを経て最終的に最も結果の出力の指示が行われている被検対象5の組み合わせを、探索すべき組み合わせとして特定する。これにより、組み合わせバンディット問題の解を自動的に探索して求めることが可能となり、情報量が増大の一途を辿る昨今において、大量の情報から高速かつ効率的に、組み合わせバンディット問題の解を求めることが可能となる。

## 【 0 0 6 2 】

なお、本発明は、上述した実施の形態に限定されるものではない。例えば、図3に示すように、2以上の解探索システム1を用いて、被検対象5を探索するようにしてもよい。図3では、2つの解探索システム1 a、解探索システム1 bを用いて被検対象5 a ~ 5 dを探索する場合について示している。ちなみに、この2つの解探索システム1 a、解探索システム1 bは、互いの同一の被検対象5 a ~ 5 dのグループの探索を行う。換言すれば図中において解探索システム1 b側の点線で示されている被検対象5 a ~ 5 dは、解探索システム1 a側の実線で示されている被検対象5 a ~ 5 dと同一のものである。

## 【 0 0 6 3 】

このケースでは、被検対象5の数は、2以上であればよい。つまり最低2つの被検対象5のうち、最良の結果の出力が期待される1の被検対象5を選択するものであってもよい。かかる点において、このケースでは、2以上の被検対象5の組み合わせを探索する場合に限定されるものではなく、1の被検対象5を解として探索するものであってもよい。

## 【 0 0 6 4 】

またこのケースにおいて被検対象5は、予め確率分布が設定されているものに限定されるものではなく、確率分布が時系列的に、つまり時間の経過に応じて変化するものであってもよい。しかも、その被検対象5における時系列的な確率分布の変化が、出力指示部4 a ~ 4 dによる出力指示に対応するものであってもよい。即ち、出力指示部4 bからの出力指示があった場合、これに対応する被検対象5 bのみが、当該出力指示に応じて自らの確率分布を変化させ、その変化させた確率分布に基づいて結果を出力する。

## 【 0 0 6 5 】

更に、各被検対象5による各確率分布の変化は、他の被検対象5による出力と独立では

10

20

30

40

50

なく、何らかの相関性を持たせるようにしてもよい。即ち、一の被検対象 5 a がより高い確率で当たりを出す確率分布とする場合に、他の一の被検対象 5 b はこれと相関性の高い確率分布を設定するようにしてもよく、更なる他の一の被検対象 5 b は、これと負の相関をもつ確率分布を設定する等してもよい。これらの相関はいかなるもので表現されていてもよい。

【 0 0 6 6 】

このケースでは、2つの解探索システム 1 a、解探索システム 1 b を用いて、上述と同様に解探索を行っていく。出力指示が行われる都度、これに応じて被検対象 5 の確率分布が時系列的に変化し、しかもその確率分布の変化が、複数の被検対象 5 との間で独立ではなく互いに相関性を持っている。これらの処理を繰り返し実行することにより、上述と同様に被検対象 5 の選択が収束され、これが探索解となる。

【 0 0 6 7 】

このケースでは、例えば株等の投資対象を被検対象 5 に当てはめて考えることもできる。株価は、上昇、下降が、他の株価と相関性を持つことが多々ある。また、出力指示部 4 により出力指示された場合は、株を買った場合（或いは空売りした場合）と考えることもできるが、仮に解探索システム 1 a のみならず、解探索システム 1 b が株を同じ購入した場合には、これに応じて株価が上昇する。これは、上述した出力指示に応じて被検対象 5 の確率分布が変化することと同様である。

【 0 0 6 8 】

更にこのケースでは、解探索システム 1 a、解探索システム 1 b 間が互いに連関していてもよい。つまり、解探索システム 1 a が特定の被検対象 5 a を選択した場合、解探索システム 1 b は、被検対象 5 c を選択するように互いに連関するように設定されていてもよい。これにより、例えば株等の投資対象を被検対象 5 に当てはめる場合において、一の解探索システム 1 a が、他の解探索システム 1 b との間で互いに連関して意思決定を行う場合が多い場合（例えば、A氏がT社の株を買った場合に、B氏は、T社と業務上の関連性が高いU社の株を買う場合）において、その探索精度を向上させることが可能となる。なお、この解探索システム 1 a による出力指示、解探索システム 1 b の出力指示の相関性はいかなるものであってもよい。

【 0 0 6 9 】

また解探索システム 1 a、解探索システム 1 b 間の連関は、互いの制御部 3 同士が直接的に相互作用するものであってもよい。かかる場合には、上述した解探索システム 1 a において計量変数を決める上での演算式を、他の解探索システム 1 b 側にも直接導入することで相関性を持たせるようにしてもよい。

【 0 0 7 0 】

なお、上述した実施形態では、2個の解探索システム 1 a、解探索システム 1 b を用いる場合を例にとり説明をしたが、これに限定されるものではなく2以上であればいかなる数で構成されていてもよい。

【実施例 1】

【 0 0 7 1 】

図 3 は、本発明を適用した解探索システム 1 を、コンピュータプログラムで実現した場合における実施例を示している。コンピュータプログラムで実現するため、戦績優劣比較部 2 は、これと同様の処理を実行する戦績優劣比較ステップに、制御部 3 は、これと同様の処理を実行する制御ステップに、出力指示部 4 は、これと同様の処理を実行する出力指示ステップとして具現化したプログラムを作成した。

【 0 0 7 2 】

被検対象 5 は、予め設定された確率の下で当選するか落選するかが決まる装置であり、例えばスロットマシン等を想定している。この被検対象 5 - 1 の当選確率 = 0 . 3 5、被検対象 5 - 2 の当選確率 = 0 . 4 5、被検対象 5 - 3 の当選確率 = 0 . 5 5、被検対象 5 - 4 の当選確率 = 0 . 6 5 とされている。また 5 - 0 は、 $x_i$  である。

【 0 0 7 3 】

10

20

30

40

50

このうち、当選確率の高い2つの被検対象5の組み合わせ（即ち、被検対象5 - 3、5 - 4）を探索解として求めることを行う。

【0074】

図3の横軸  $t$  はフィードバック回数である。また、戦績指数  $q_i$  は、(1)式に基づいて被検対象5毎に求めている。優劣  $S_i$  は、(2)式に基づいて被検対象5毎に求めている。  $a_i$  は、表2から導き出される - 1、0、1の何れかの数値である。

【0075】

また  $v_i$  は、  $a_i$  に基づいて計量変数を増やす方向に推進させるのか、或いは減らす方向に推進させるのかを決めるパラメータである。仮に  $a_i$  が計量変数を増減させる際の加速度であるとすれば、この  $v_i$  は、計量変数の増減速度である。  $a_i$  は、  $v_i$  の1回微分で表すことができる。更に、計量変数  $x_i$  は、実際に計量変数の増減速度  $v_i$  に基づいて増減させられた実際のパラメータ値である。ここで閾値は0としている。

10

【0076】

このようなプログラムによる結果より、フィードバック回数が増加するにつれて、計量変数  $x_i$  は、殆ど被検対象5 - 3、5 - 4のみが大きくなってそのまま出力指示部4によって出力指示が行われることが繰り返されるのがわかる。

【0077】

このように本発明を適用したコンピュータプログラムにより、組み合わせバンディット問題を自動的に求めることができることを検証することができた。

【0078】

20

図4は、他の実施例ではあるが、それぞれ被検対象5の当選確率が、0.2、0.5、0.8の3つのサンプルについて、本発明を適用したコンピュータプログラムにより、より当選確率の高い組み合わせを選択するシミュレーションを行った結果である。実線が本発明に相当する。比較例1は、ソフトマックスアルゴリズムを組み合わせ報酬最大化問題に拡張したものである。また比較例2は、イプシロングリーディアルゴリズムを組み合わせ報酬最大化問題に拡張したものである。その結果、本発明を適用したコンピュータプログラムは、フィードバック回数が増加するにつれて、他の比較例よりも正答率が高くなるのが分かる。

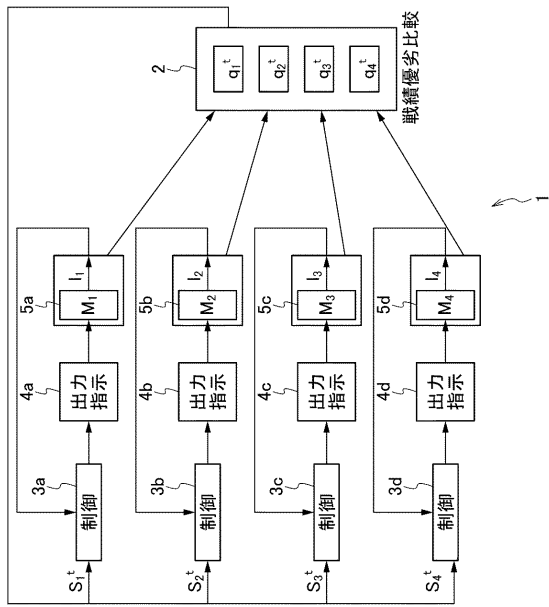
【符号の説明】

【0079】

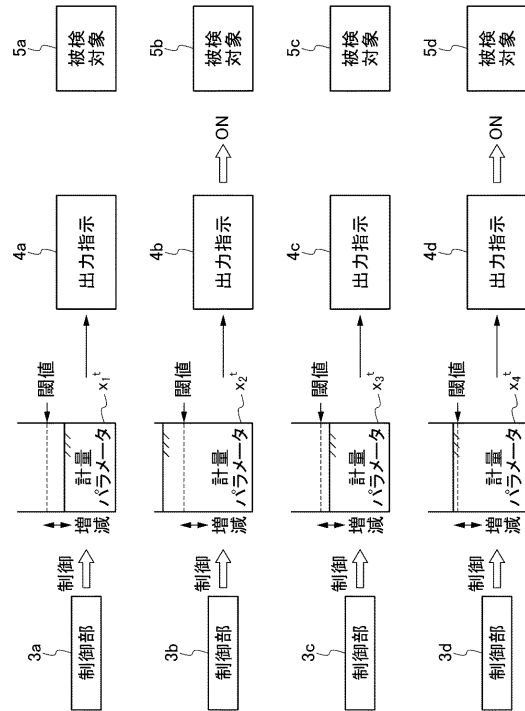
30

- 1 解探索システム
- 2 戦績優劣比較部
- 3 制御部
- 4 出力指示部
- 5 被検対象

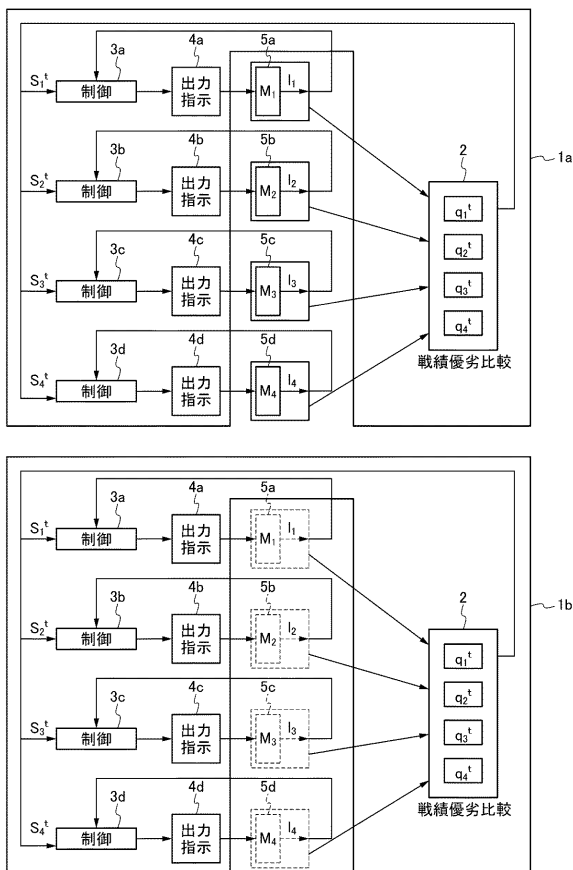
【図1】



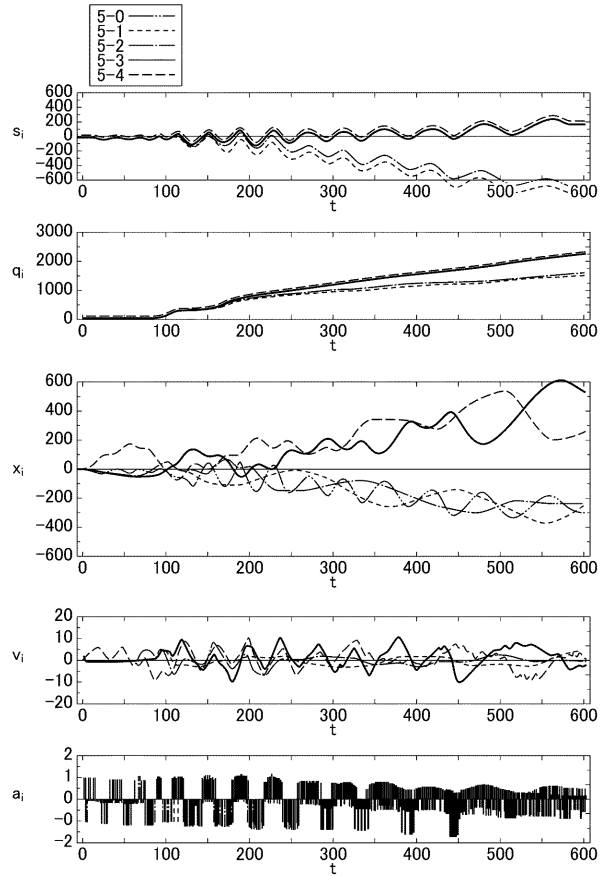
【図2】



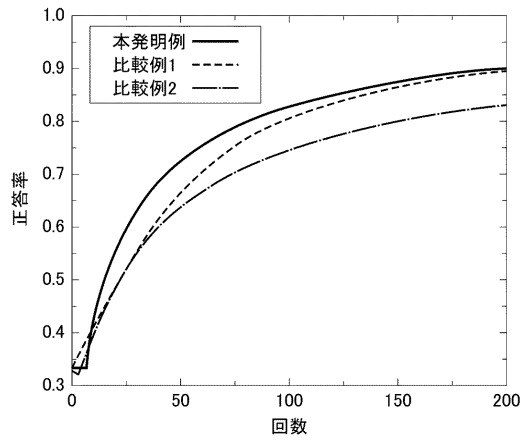
【図3】



【図4】



【図 5】



---

フロントページの続き

(72)発明者 原 正彦

埼玉県和光市広沢 2 番 1 号 独立行政法人理化学研究所内

審査官 松野 広一

(56)参考文献 米国特許出願公開第 2 0 0 8 / 0 1 4 0 5 9 1 ( U S , A 1 )

AUER, Peter et al, Finite time Analysis of the Multiarmed Bandit Problem, Machine learning, 2 0 0 2 年, Vol.47, pp.235 256

VERMOREL, Joanne s, MOHRI, Mehryar, Multi Armed Bandit Algorithms and Empirical Evaluation, In Proceedings of the 16th European Conference on Machine Learning, 2 0 0 5 年 1 0 月, Vol.3720, pp.437 448, U R L , <http://www.cs.nyu.edu/~mohri/pub/bandit.pdf>

大用 庫智 外 2 名, 非定常 N 本腕バンディット問題に対する人間の認知バイアスの適用, 2 0 1 1 年度人工知能学会全国大会 ( 第 2 5 回 ) 論文集 [ C D - R O M ] 2 0 1 1 年度人工知能学会全国大会, 日本, 2 0 1 1 年 6 月 1 日, pp.1 4

青野 真士 外 4 名, 粘菌コンピュータと確率探索アルゴリズム, システム / 制御 / 情報, 日本, システム制御情報学会, 2 0 1 1 年 1 2 月 1 5 日, Vol.55 No.12, pp.26 31

(58)調査した分野(Int.Cl., D B 名)

G 0 6 Q 1 0 / 0 0 - 9 9 / 0 0

G 0 6 F 1 9 / 0 0

G 0 6 N 3 / 0 0

G 0 6 N 5 / 0 4