

(12) 特許協力条約に基づいて公開された国際出願

(19) 世界知的所有権機関
国際事務局

(43) 国際公開日
2020年6月18日(18.06.2020)



(10) 国際公開番号
WO 2020/121494 A1

- (51) 国際特許分類:
G06N 20/00 (2019.01)
- (21) 国際出願番号: PCT/JP2018/045947
- (22) 国際出願日: 2018年12月13日(13.12.2018)
- (25) 国際出願の言語: 日本語
- (26) 国際公開の言語: 日本語
- (71) 出願人: 日本電気株式会社 (NEC CORPORATION) [JP/JP]; 〒1088001 東京都港区芝五丁目7番1号 Tokyo (JP).
- (72) 発明者: 森 達哉(MORI Tatsuya); 〒1088001 東京都港区芝五丁目7番1号 日本電気株式会社内 Tokyo (JP). 平岡 拓也(HIRAOKA Takuya); 〒1088001 東京都港区芝五丁目7番1号 日

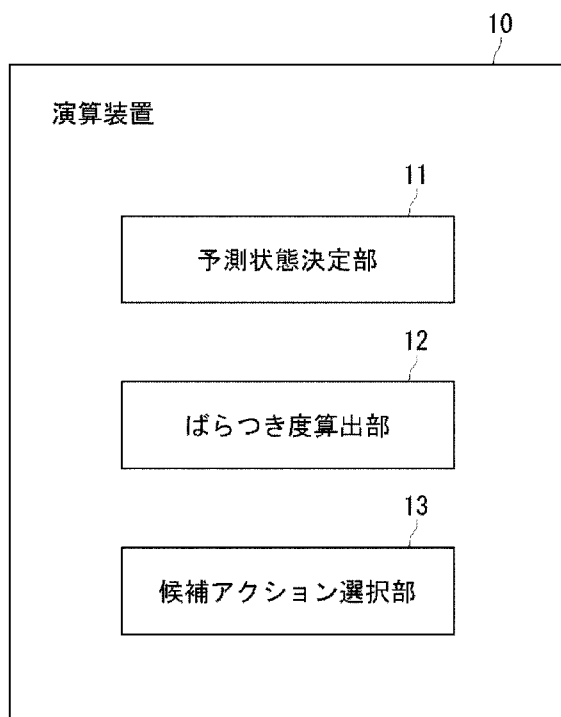
本電気株式会社内 Tokyo (JP). タンカラットブット(TANGKARATT Voot); 〒3510198 埼玉県和光市広沢2番1号 国立研究開発法人理化学研究所内 Saitama (JP).

(74) 代理人: 家入 健(IEIRI Takeshi); 〒2210835 神奈川県横浜市神奈川区鶴屋町三丁目3番8 アサヒビルディング5階 響国際特許事務所 Kanagawa (JP).

(81) 指定国(表示のない限り、全ての種類の国内保護が可能): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DJ, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JO, JP, KE, KG, KH,

(54) Title: ARITHMETIC DEVICE, ACTION DETERMINATION METHOD, AND NON-TRANSITORY COMPUTER-READABLE MEDIUM STORING CONTROL PROGRAM

(54) 発明の名称: 演算装置、アクション決定方法、及び制御プログラムを格納する非一時的なコンピュータ可読媒体



- 10 Arithmetic device
11 Predictive state determination unit
12 Dispersion degree calculation unit
13 Candidate action selection unit

(57) Abstract: In an arithmetic device (10), a predictive state determination unit (11) determines, using a plurality of transition information units, a plurality of predictive states pertaining to each of a plurality of candidate actions possible in a first state. A dispersion degree calculation unit (12) calculates the degree of dispersion of the plurality of predictive states

KN, KP, KR, KW, KZ, LA, LC, LK, LR, LS, LU, LY,
MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ,
NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT,
QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL,
SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA,
UG, US, UZ, VC, VN, ZA, ZM, ZW.

- (84) 指定国(表示のない限り、全ての種類の広域保護が可能): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), ユーラシア (AM, AZ, BY, KG, KZ, RU, TJ, TM), ヨーロッパ (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

添付公開書類 :

- 一 国際調査報告 (条約第21条(3))

determined pertaining to each candidate action by the predictive state determination unit (11). A candidate action selection unit (13) selects some of the plurality of candidate actions on the basis of multiple degrees of dispersion calculated by the dispersion degree calculation unit (12).

(57) 要約 : 演算装置 (10) にて予測状態決定部 (11) は、複数の遷移情報ユニットを用いて、第1状態にて可能な複数の候補アクションのそれぞれに関して複数の予測状態を決定する。ばらつき度算出部 (12) は、予測状態決定部 (11) にて各候補アクションに関して決定された複数の予測状態のばらつき度を算出する。候補アクション選択部 (13) は、ばらつき度算出部 (12) にて算出された複数のばらつき度に基づいて、上記の複数の候補アクションのうちの一部の候補アクションを選択する。

明 細 書

発明の名称：

演算装置、アクション決定方法、及び制御プログラムを格納する非一時的なコンピュータ可読媒体

技術分野

[0001] 本開示は、演算装置、アクション決定方法、及び制御プログラムに関する。

背景技術

[0002] 「強化学習」に関して種々の研究が行われている（例えば、非特許文献1）。強化学習の目的の1つは、実環境に対して時系列的に複数の行動（アクション）を行った結果、実環境から得られる「累積報酬」を最大化する、方策（Policy）を学習することである。

先行技術文献

非特許文献

[0003] 非特許文献1：Richard S. Sutton and Andrew G. Barto, “Reinforcement Learning : An Introduction”, Second Edition, MIT Press, 2018

発明の概要

発明が解決しようとする課題

[0004] ところで、好適な方策を効率良く学習するためには、実環境の状態についての「状態空間」を効率的に探索する必要がある。

[0005] しかしながら、非特許文献1では探索の重要性について言及されているが、効率的な探索を実現する具体的な技術については開示されていない。

[0006] 本開示の目的は、効率的な探索を実現できる、演算装置、アクション決定方法、及び制御プログラムを提供することにある。

課題を解決するための手段

- [0007] 第1の態様にかかる演算装置は、第1タイミングでの第1状態と、前記第1タイミング以降の第2タイミングでの第2状態との関係性を表す遷移情報を複数用いて、第1状態にて可能な複数の候補アクションのそれぞれに関して複数の第2状態を決定する決定手段と、各前記候補アクションに関して、前記複数の第2状態のばらつき度を算出する算出手段と、前記ばらつき度に基づき、前記複数の候補アクションから一部の候補アクションを選択する選択手段と、を具備する。
- [0008] 第2の態様にかかるアクション決定方法は、情報処理装置によって、第1タイミングでの第1状態と、前記第1タイミング以降の第2タイミングでの第2状態との関係性を表す遷移情報を複数用いて、第1状態にて可能な複数の候補アクションのそれぞれに関して複数の第2状態を決定し、各前記候補アクションに関して、前記複数の第2状態のばらつき度を算出し、前記ばらつき度に基づき、前記複数の候補アクションから一部の候補アクションを選択する。
- [0009] 第3の態様にかかる制御プログラムは、第1タイミングでの第1状態と、前記第1タイミング以降の第2タイミングでの第2状態との関係性を表す遷移情報を複数用いて、第1状態にて可能な複数の候補アクションのそれぞれに関して複数の第2状態を決定し、各前記候補アクションに関して、前記複数の第2状態のばらつき度を算出し、前記ばらつき度に基づき、前記複数の候補アクションから一部の候補アクションを選択する処理を、演算装置に実行させる。

発明の効果

- [0010] 本開示により、効率的な探索を実現できる、演算装置、アクション決定方法、及び制御プログラムを提供することができる。

図面の簡単な説明

- [0011] [図1]第1実施形態の演算装置の一例を示すブロック図である。
[図2]第2実施形態の演算装置を含む制御装置の一例を示すブロック図である

- 。
- [図3]第2実施形態の演算装置の処理動作の一例を示すフローチャートである
- 。
- [図4]第3実施形態の演算装置を含む制御装置の一例を示すブロック図である
- 。
- [図5]第3実施形態の演算装置の処理動作の一例を示すフローチャートである
- 。
- [図6]演算装置のハードウェア構成例を示す図である。

発明を実施するための形態

- [0012] 以下、図面を参照しつつ、実施形態について説明する。なお、実施形態において、同一又は同等の要素には、同一の符号を付し、重複する説明は省略される。
- [0013] <第1実施形態>
- 図1は、第1実施形態の演算装置の一例を示すブロック図である。図1において演算装置（アクション決定装置）10は、予測状態決定部11と、ばらつき度算出部12と、候補アクション選択部13とを有している。
- [0014] 説明の便宜上、あるタイミング（以降、「第1タイミング」と表す）における制御対象の状態を「第1状態」と表す。あるタイミング以降の1つタイミング（以降、「第2タイミング」と表す）における制御対象の状態を「第2状態」と表す。制御対象の状態は、第1状態に応じたアクションが実施された後に第2状態に変化するとする。また、第1状態と、第2状態とは、必ずしも、相互に異なる状態である必要はなく、同じ状態を表していてもよい。以降の説明においては、説明の便宜上、第1状態と、第2状態との異同によらずに、「第1状態から第2状態に変化する」と表すこととする。また、第1タイミング、及び、第2タイミングは、特定のタイミングを表しているわけではなく、相互に異なる2つのタイミングを表している。
- [0015] 予測状態決定部11は、複数の状態遷移情報（遷移情報ユニット）を用いて、第1状態にて可能な複数の「候補アクション」のそれぞれに関して複数

の「予測状態」を決定する。各遷移情報ユニットは、第1状態と、該第1状態におけるアクションとから、該第1タイミング以降（たとえば、第2タイミング）の予測状態を算出するために用いられる。すなわち、各遷移情報ユニットは、各遷移情報ユニットの第1状態を保持しており、該第1状態及びアクションの組み合わせに応じた予測状態を決定する機能を有している。ここで、例えば、各遷移情報ユニットは、あるタイミングでの実環境の状態（実環境状態）と、該あるタイミングで実環境に対して実際に行われたアクションとが関連付けされたセットを含む「履歴情報」に基づいて作成（訓練）される。該セットは、2つの状態と、当該2つの状態間におけるアクションとが関連付けされた情報を表す。

[0016] ばらつき度算出部12は、予測状態決定部11にて各候補アクションに関して決定された複数の予測状態の「ばらつき度」を算出する。ここでは、第1状態にて可能な候補アクションは複数存在しているので、複数の候補アクションにそれぞれ対応する複数のばらつき度が算出されることになる。「ばらつき度」は、例えば、分散値である。

[0017] 候補アクション選択部13は、ばらつき度算出部12にて算出された複数のばらつき度に基づいて、上記の複数の候補アクションのうちの一部の候補アクションを選択する。例えば、候補アクション選択部13は、上記の複数の候補アクションのうちで、ばらつき度算出部12にて算出された複数のばらつき度のうちの最大値に対応する候補アクションを選択する。

[0018] 以上のように第1実施形態によれば、演算装置10にて予測状態決定部11は、複数の遷移情報ユニットを用いて、第1状態にて可能な複数の「候補アクション」のそれぞれに関して複数の「予測状態」を決定する。ばらつき度算出部12は、予測状態決定部11にて各候補アクションに関して決定された複数の予測状態の「ばらつき度」を算出する。候補アクション選択部13は、ばらつき度算出部12にて算出された複数のばらつき度に基づいて、上記の複数の候補アクションのうちの一部の候補アクションを選択する。

[0019] この演算装置10の構成により、効率的な探索を行うことができる。すな

わち、候補アクションによる第1状態から第2状態への状態遷移が遷移情報ユニットにおいて「訓練の不十分な状態遷移」である場合、その状態遷移の予測状態についての「ばらつき度」が高くなる傾向にある。すなわち、「ばらつき度」は、遷移情報ユニットにおける状態遷移の訓練進捗度を示す指標として用いることができる。また、上記「訓練の不十分な状態遷移」は、上記「履歴情報」に十分に蓄積されていない状態遷移、つまり、実環境において探索が十分でない状態遷移を表していることもある。このため、ばらつき度に基づき候補アクションを選択することによって、探索が十分でない状態遷移（つまり、状態及びアクションの組み合わせ）を積極的に探索することができる。よって、効率的に探索を行なうことができる。さらに、探索が十分でない状態遷移を積極的に探索することができるので、遷移情報ユニットの訓練を効率的に行うことができる。

[0020] <第2実施形態>

第2実施形態は、より具体的な実施形態に関する。

[0021] <制御装置の概要>

図2は、第2実施形態の演算装置30を含む制御装置20の一例を示すブロック図である。図2には、制御装置20の他に、指令実行装置50及び制御対象60が図示されている。

[0022] たとえば、制御対象60が車両である場合に、制御装置20は、たとえば、エンジンの回転数や、車両の速度や、周囲の状況等の観測値（特徴量）に基づき、ハンドルを右に回す、アクセルを踏む、ブレーキを踏む等のアクションを決定する。指令実行装置50は、演算装置30によって決定されたアクションに従いアクセル、ハンドル、または、ブレーキを制御する。

[0023] たとえば、制御対象60が発電機である場合に、制御装置20は、たとえば、タービンの回転数や、燃焼炉の温度や、燃焼炉の圧力等の観測値に基づき、燃料の量を増やす、燃料の量を減らす等のアクションを決定する。指令実行装置50は、制御装置20によって決定されたアクションに従い燃料の量を調整するバルブを閉める、あるいは、バルブを開く等の制御を実行する

- 。
- [0024] 制御対象60は、上述した例に限定されず、たとえば、生産工場や、化学工場であってもよいし、車両の動作や、発電機における動作などをシミュレーションしているシミュレータ等であってもよい。
- [0025] 観測値に基づきアクションを決定する処理については、図3を参照しながら後述する。
- [0026] 制御装置20は、後述するような、「処理フェーズ1」、「処理フェーズ2」、及び「処理フェーズ3」を実行する。制御装置20は、これらの処理を実行することによって、制御対象60の状態 (state) が、より早期に所望状態に近づくようアクションを決定する。この際に、制御装置20は、制御対象60の状態に対して実行するアクションを、方策 (Policy) 情報と、報酬 (reward) 情報とに基づき決定する。
- [0027] 方策情報は、制御対象60がある状態である場合に、実行可能なアクションを表す。方策情報は、たとえば、該ある状態と、該アクションとが関連付けされた情報を用いて実現することができる。方策情報は、たとえば、該ある状態を与えた場合に該アクションを算出する処理であってもよい。当該処理は、たとえば、ある関数、または、統計的な手法によって算出された、当該ある状態と、当該アクションとの関係性を表すモデルであってもよい。すなわち、方策情報は、上述した例に限定されない。
- [0028] 報酬情報は、ある状態が望ましい程度 (以降、「報酬程度」と表す) を表す。報酬情報は、たとえば、該ある状態と、該程度とが関連付けされた情報を用いて実現することができる。報酬情報は、たとえば、該ある状態を与えた場合に該報酬程度を算出する処理であってもよい。当該処理は、たとえば、ある関数、または、統計的な手法によって算出された、当該ある状態と、当該報酬程度との関係性を表すモデルであってもよい。すなわち、報酬情報は、上述した例に限定されない。
- [0029] 以降の説明においては、説明の便宜上、制御対象60は、車両や、発電機等 (以降、「実環境」と表す) であるとする。あるタイミング (以降、「第

1 タイミング」と表す)における制御対象60の状態を「第1状態」と表す。あるタイミングの次のタイミング(以降、「第2タイミング」と表す)における制御対象60の状態を「第2状態」と表す。制御対象60の状態は、第1状態に応じたアクションが実施された後に第2状態に変化するとする。また、第1状態と、第2状態とは、必ずしも、相互に異なる状態である必要はなく、同じ状態を表していてもよい。以降の説明においては、説明の便宜上、第1状態と、第2状態との異同によらずに、「第1状態から第2状態に変化する」と表すこととする。

[0030] 制御装置20は、複数のタイミングに関して、制御対象60の観測値を参照しながら、処理フェーズ1乃至処理フェーズ3にて後述するような処理を実行することによって、タイミングごとにアクションを決定する。すなわち、制御装置20は、第1タイミングに関して処理を実行した後に、第2タイミングに関して処理を実行し、さらに、第2タイミングより後のタイミングに関して処理を実行する。したがって、第1タイミング、及び、第2タイミングは、特定のタイミングを表しているわけではなく、制御装置20における処理に関して連続している2つのタイミングを表している。

[0031] (処理フェーズ1)

制御装置20は、状態遷移情報(後述する)に基づき第1状態である制御対象60に関して、アクションを実行した後の制御対象60の第2状態を推定する。制御装置20は、複数の候補アクションに関して、それぞれ、第2状態を推定する処理を実行する。その後、制御装置20は、報酬情報を用いて、推定した各第2状態に関する報酬程度を算出する。制御装置20は、複数の候補アクションの中から、算出した報酬程度が上位の候補アクションのうちの一つアクションを選択する。制御装置20は、複数の候補アクションの中から、算出した報酬程度が最も大きなアクションを一つ選択してもよい。制御装置20は、選択したアクションを示す制御指令を、指令実行装置50へ出力する。

[0032] 上位は、たとえば、報酬程度が最も大きいものから、報酬程度が大きい順

に数えて、1%、5%、または、10%等の所定の割合以内であることを表している。

[0033] ここで、状態遷移情報について説明する。状態遷移情報は、第1状態と、第2状態との間の関係性を表す情報である。状態遷移情報は、第1状態と、第2状態とが関連付けされた情報であってもよいし、第1状態と、第2状態とが関連付けされた訓練データを用いたニューラルネットワーク等の統計的な手法によって算出された情報であってもよい。状態遷移情報は、さらに、第1状態にて実行可能なアクションを表す情報を含んでいてもよく、上述した例に限定されない。

[0034] 指令実行装置50は、制御装置20によって制御指令を受け取り、受け取った該制御指令が示すアクションを、制御対象60に関して実行する。この結果、制御対象60の状態は、第1状態から第2状態に変化する。

[0035] 説明の便宜上、制御対象60には、制御対象60を観測しているセンサー（図示せず）が取り付けられているとする。センサーは、制御対象60に関して観測した観測値を表すセンサー情報を作成し、作成したセンサー情報を出力するとする。制御対象60を観測しているセンサーは、複数であってもよい。

[0036] 制御装置20は、第1状態に関するアクションが実行された後に、センサーによって作成された該センサー情報を受け取り、受け取った該センサー情報に関する第2状態を決定する。制御装置20は、該第1状態と、該アクションと、該第2状態とが関連付けされた情報（以降、「履歴情報」と表す）を作成する。制御装置20は、作成した履歴情報を、後述する履歴情報記憶部41に格納してもよい。

[0037] 処理フェーズ1に関して上述したような処理が、複数のタイミングに関して実行されることにより、後述する履歴情報記憶部41には、複数のタイミングにおける履歴情報が蓄積される。

[0038] （処理フェーズ2）

制御装置20は、処理フェーズ1にて蓄積された履歴情報を用いて、状態

遷移情報を更新する（または、作成する）。状態遷移情報を、ニューラルネットワークを用いて作成する場合に、制御装置20は、上述したような履歴情報に含まれているデータを訓練データとして用いて、当該状態遷移情報を作成する。後述するように、制御装置20は、たとえば、構成が相互に異なっているニューラルネットワークを用いて、複数の状態遷移情報を作成する。

[0039] （処理フェーズ3）

制御装置20は、複数の候補アクションについて、候補アクションをそれぞれ対象に関して施した後における第2状態を、状態遷移情報に基づき予測する。制御装置20は、相互に異なる状態遷移情報（すなわち、各遷移情報ユニット）を用いることによって、複数の第2状態を予測する。説明の便宜上、第2状態と、予測された第2状態とを区別するため、予測された第2状態を「擬似状態」と表す。すなわち、制御装置20は、相互に異なる状態遷移情報（すなわち、各遷移情報ユニット）を用いることによって、擬似状態を作成する。

[0040] 状態遷移情報を、ニューラルネットワークを用いて作成する場合に、制御装置20は、第1状態、及び、当該第1状態における候補アクションを表す情報のうち、少なくともいずれかに対して当該状態遷移情報を適用することによって、該擬似状態を作成する。

[0041] 処理フェーズ3に関して上述した処理によって、制御装置20は、各候補アクションに関して、複数の擬似状態を作成する。制御装置20は、各候補アクションに関して、複数の擬似状態のばらつき度を算出する。

[0042] 制御装置20は、複数の候補アクションの中から、該ばらつき度に基づきアクションを選択する。制御装置20は、複数の候補アクションの中から、算出したばらつき度が上位である候補アクションを特定し、特定した候補アクションの中からアクションを選択する。制御装置20は、たとえば、複数の候補アクションの中から、算出したばらつき度が最も大きな候補アクションを選択してもよい。

[0043] 上位は、たとえば、ばらつき度が最も大きいものから、ばらつき度が大きい順に数えて、1%、5%、または、10%等の所定の割合以内であることを表している。

[0044] 制御装置20は、報酬情報を用いて、1つのアクション後の擬似状態における報酬程度を求め、求めた報酬程度と、当該1つのアクションに対するばらつき度とに基づき、アクションを選択してもよい。

[0045] 擬似状態が複数である場合に、制御装置20は、たとえば、各擬似状態に関する報酬程度の平均（または、中央値）を求めることによって、アクションに関する報酬程度を求める。または、制御装置20は、たとえば、各擬似状態の頻度が上位の状態を求め、求めた状態に関する報酬程度の平均（または、中央値）を求めることによって、アクションに関する報酬程度を求める。この場合に、上位は、たとえば、頻度が最も高いものから、頻度が高い順に数えて、1%、5%、または、10%等の所定の割合以内であることを表している。アクションに関する報酬程度を求める処理は、上述した例に限定されない。

[0046] また、1つのアクションに関する報酬程度と、該1つのアクションに関するばらつき度とに基づき、アクションを選択する処理は、たとえば、該報酬程度と、該ばらつき度とを足し算してもよいし、該報酬程度と、該ばらつき度との重み付き平均を算出してもよい。アクションを選択する処理は、上述した例に限定されない。

[0047] 制御装置20は、アクションを選択した後に、選択したアクションを示す制御指令を指令実行装置50へ出力する。指令実行装置50は、受け取った制御指令が示すアクションを制御対象60に関して実行する。

[0048] <制御装置の構成例>

図2において制御装置20は、演算装置30と、記憶装置40とを有している。演算装置30は、状態推定部31と、状態遷移情報更新部（状態遷移情報作成部）32と、制御指令演算部33と、予測状態決定部11と、ばらつき度算出部12と、候補アクション選択部13とを有している。記憶装置

40は、履歴情報記憶部41と、状態遷移情報記憶部42と、方策情報記憶部43とを有している。

[0049] (処理フェーズ1)

状態推定部31は、制御対象60の第1状態を表す観測値(パラメタ値、センサー情報)を受け取る。状態推定部31は、受け取ったセンサー情報と、状態遷移情報とに基づき、第1状態である制御対象60に関してアクションを実行した後の制御対象60の第2状態を推定する。状態推定部31は、複数の候補アクションにおけるアクションに関して、それぞれ、第2状態を推定する処理を実行する。すなわち、状態推定部31は、各候補アクションに関して擬似状態を作成する。

[0050] 制御指令演算部33は、報酬情報を用いて、状態推定部31によって作成された各擬似状態に関する報酬程度を算出する。制御指令演算部33は、複数の候補アクションの中から、算出した報酬程度が上位の候補アクションのうちの1つアクションを選択する。制御指令演算部33は、選択したアクションを示す制御指令を作成し、作成した制御指令を指令実行装置50へ出力する。

[0051] 指令実行装置50は、制御指令を受け取り、受け取った制御指令が示すアクションに従い、制御対象60に関するアクションを実行する。制御対象60に関するアクションの結果、制御対象60の状態は、第1状態から第2状態に変化する。

[0052] 状態推定部31は、制御対象60の状態(この場合、第2状態)を表す観測値(パラメタ値、センサー情報)を受け取る。状態推定部31は、第1状態と、第1状態にて実行されたアクションと、該第2状態とが関連付けされた履歴情報を作成し、作成した履歴情報を履歴情報記憶部41に格納する。

[0053] 処理フェーズ1に関して上述したような処理を繰り返すことによって、履歴情報記憶部41には、上記の履歴情報が蓄積される。

[0054] (処理フェーズ2)

説明の便宜上、ニューラルネットワーク等の統計的な手法(所定の処理手

順)を用いて状態遷移情報を作成する例を用いて、処理フェーズ2における処理を説明する。所定の処理手順は、例えば、ニューラルネット等の機械学習法に従った手順である。

[0055] 状態遷移情報更新部32は、履歴情報記憶部41に蓄積されている履歴情報を用いて、所定の処理手順に従って、複数の遷移情報ユニットを作成する。すなわち、状態遷移情報更新部32は、該履歴情報を訓練データとして、所定の処理手順に従い状態遷移情報を作成し、作成した状態遷移情報を状態遷移情報記憶部42に格納する。上述したように、状態遷移情報は、第1状態と、第2状態との関係性を表す。

[0056] 例えば、状態遷移情報更新部32は、構成が互いに異なる複数のニューラルネットを用いて、複数の遷移情報ユニットを作成してもよい。構成が互いに異なる複数のニューラルネットは、例えば、互いにノードの数又はノード間の接続パターンが異なる複数のニューラルネットである。また、互いに構成の異なる複数のニューラルネットは、あるニューラルネットワークと、当該あるニューラルネットワークにおける一部のノードが存在していない(すなわち、一部のノードがドロップアウトしている)ニューラルネットワークとを用いて実現されていてもよい。

[0057] 状態遷移情報更新部32は、パラメタの初期値が異なる複数のニューラルネットを用いて、複数の遷移情報ユニットを作成してもよい。

[0058] 状態遷移情報更新部32は、履歴情報のうちの一部のデータ、または、履歴情報から重複を許してサンプリングしたものを訓練データとして用いてもよい。この場合に、複数の遷移情報ユニットは、相互に異なる訓練データに対して状態遷移情報を作成する。

[0059] なお、所定の処理手順は、ニューラルネットに限定されない。例えば所定の処理手順は、SVM (support vector machine)、ランダムフォレスト、バギング (bootstrap aggregating)、又は、ベイジアンネットワークを算出する手順であってもよい。

[0060] (処理フェーズ3)

予測状態決定部 11 は、複数の候補アクションについて、候補アクションをそれぞれ対象に関して施した後における第 2 状態を、状態遷移情報に基づき予測する。予測状態決定部 11 は、相互に異なる状態遷移情報（すなわち、各遷移情報ユニット）を用いることによって、複数の疑似状態を作成する。

[0061] ばらつき度算出部 12 は、予測状態決定部 11 によって作成された複数の疑似状態のばらつき度（たとえば、分散値、エントロピー等）を算出し、算出したばらつき度を候補アクション選択部 13 へ出力する。ばらつき度は、たとえば、分散値にある数を加えた値等であってもよく、上述した例に限定されない。

[0062] 候補アクション選択部 13 は、複数の候補アクションの中から、該ばらつき度に基づきアクションを選択する。候補アクション選択部 13 は、複数の候補アクションの中から、算出したばらつき度が上位である候補アクションを特定し、特定した候補アクションの中からアクションを選択する。候補アクション選択部 13 は、たとえば、複数の候補アクションの中から、算出したばらつき度が最も大きな候補アクションを選択してもよい。

[0063] 制御指令演算部 33 は、候補アクション選択部 13 が選択したアクションを示す制御指令を作成し、作成した制御指令を指令実行装置 50 へ出力する。

[0064] 上述したように候補アクション選択部 13 は、ばらつき度が大きいアクションを選択する。ばらつき度は、状態遷移情報に従い算出された結果がばらついていることを表している。このため、ばらつき度が大きい場合には、状態遷移情報が不安定であることを表しているということもできる。すなわち、ばらつき度が大きいアクションを実行することによって、探索が十分でない状態遷移を積極的に探索することができるという効果を奏する。

[0065] 候補アクション選択部 13 は、状態価値情報に基づき、状態に関する価値の程度を表す状態価値情報を作成してもよい。状態価値情報は、たとえば、状態に対して、当該状態の価値の程度を表す関数である。この場合に、価値

は、当該状態を実現することが望ましい程度を表す情報であるとも言えることができる。状態価値情報は、アクション後における制御対象60の状態がどの程度望ましいのかを表す情報ともいうことができる。状態価値情報は、また、当該アクションがどの程度の望ましいのかを表す情報ともいうことができる。

[0066] 候補アクション選択部13は、状態価値情報を作成する処理において、報酬情報を用いてもよい。たとえば、候補アクション選択部13は、各アクションに関して算出されたばらつき度を、新たに、状態価値情報として設定してもよい。たとえば、候補アクション選択部13は、各アクションに関して算出されたばらつき度を状態価値情報として設定し、その後、当該アクションに関する報酬情報を加える等の処理を実行することによって、状態価値情報を更新してもよい。この場合に、ばらつき度は、報酬情報に対する追加的な報酬（疑似追加報酬）であるともいうことができる。

[0067] 状態価値情報を作成する処理は、上述した例に限定されず、たとえば、報酬情報に所定の値を加算した値、報酬情報に所定の値を減算した値、または、報酬情報に所定の値を乗算した値等に基づき実行されてもよい。すなわち、ばらつき度が大きいほど、状態価値情報は、価値の程度が高いことを表す情報であればよい。

[0068] 候補アクション選択部13は、状態価値情報に基づき、複数の候補アクションの中から、該価値の程度が上位の候補アクションを選択し、選択した候補アクションからアクションを選択してもよい。候補アクション選択部13は、たとえば、算出した価値が最も高い候補アクションを選択してもよい。この場合に、上位は、たとえば、価値の程度が最も高いものから価値の程度が高い順に数えて、1%、5%、または、10%等の所定の割合以内であることを表している。

[0069] 制御指令が作成された後に、指令実行装置50は、該制御指令を受け取り、受け取った制御指令が示すアクションに従い、制御対象60に関するアクションを実行する。制御対象60に関するアクションの結果、制御対象60

の状態は、第1状態から第2状態に変化する。

[0070] 状態推定部31は、制御対象60の状態（この場合、第2状態）を表す観測値（パラメタ値、センサー情報）を受け取る。状態推定部31は、第1状態と、第1状態にて実行されたアクションと、該第2状態とが関連付けされた履歴情報を作成し、作成した履歴情報を履歴情報記憶部41に格納する。

[0071] 処理フェーズ3に関して上述したような処理が、複数のタイミングに関して実行されることにより、履歴情報記憶部（不図示）には、複数のタイミングにおける履歴情報が蓄積される。

[0072] <制御装置の動作例>

以上の構成を有する演算装置30の処理動作の一例について説明する。図3は、第2実施形態の演算装置の処理動作の一例を示すフローチャートである。図3に示すフローチャートにおいて、ステップS101は、上記の処理フェーズ1に対応し、ステップS102は、処理フェーズ2に対応し、ステップS103、S104は、処理フェーズ3に対応する。

[0073] 演算装置30は、履歴情報が蓄積されるまで、処理フェーズ1及び処理フェーズ2、または、処理フェーズ3及び処理フェーズ2のうち、少なくとも、いずれかの処理を繰り返すことによって、履歴情報を取得する（ステップS101）。

[0074] 演算装置30は、処理フェーズ2に示された処理に従い、状態遷移情報を更新する（ステップS102）。

[0075] 演算装置30は、処理フェーズ3にて上述した処理に従い、ばらつき度を算出する（ステップS103）。

[0076] 演算装置30は、履歴情報に基づき方策情報を更新する（ステップS104）。具体的には、演算装置30は、履歴情報に基づき、第1状態と、当該第1状態にて実行したアクションと、第2状態と特定し、特定したこれらの情報を用いて、方策情報を更新する。そして、処理ステップは、ステップS101（処理フェーズ1）に戻る。

[0077] なお、以上の説明では、演算装置30が、処理フェーズ3にて、履歴情報

を蓄積してから方策情報を更新し、その後直ぐに、処理フェーズ1に戻るものとして説明した。説明の便宜上、本実施形態においては、図3を参照しながら上述した処理を、「バッチ学習」と記載する。

すなわち、バッチ学習は、ある程度（説明の便宜上、「第1蓄積程度」と称する）の履歴情報が蓄積されてから、該履歴情報を用いて方策情報を更新（または、作成）する処理を表す。第1蓄積程度は、履歴が複数であることを表している。ただし、演算装置30における処理は、上述したバッチ学習に限定されず、例えば、方策情報は、オンライン学習によって更新（または、作成）されてもよい、ミニバッチ学習によって更新（または、作成）されてもよい。

[0078] オンライン学習は、履歴情報に履歴が1つ追加されるごとに、該履歴情報を用いて方策情報を更新（または、作成）する処理を表す。

[0079] ミニバッチ学習は、ある程度（説明の便宜上、「第2蓄積程度」と称する）の履歴情報が蓄積されてから、該履歴情報を用いて方策情報を更新（または、作成）する処理を表す。第2蓄積程度は、履歴が複数であることを表している。ミニバッチ学習は、バッチ学習と類似した処理である。しかし、第2蓄積程度は、第1蓄積程度に比べて少ない。

[0080] 第1蓄積程度、及び、第2蓄積程度は、必ずしも、処理フェーズ1乃至処理フェーズ3に示された反復処理ごとに一定の程度でなくともよく、該反復処理ごと異なる個数を表していてもよい。

[0081] オンライン学習の場合、履歴情報を取得する度に方策情報を更新して、ステップS101（処理フェーズ1）へ戻るように、修正されてもよい。すなわち、オンライン学習の場合には、候補アクション選択部13は、第2状態に関するセンサー情報が届く度に、ポリシーモデルを更新する。

[0082] 「ミニバッチ学習」は、方策情報の更新タイミング以外は、上記「オンライン学習」の処理動作と変わらない。すなわち、「ミニバッチ学習」にて一度の方策情報の更新に用いられる履歴情報量は、「オンライン学習」よりも多いので、「ミニバッチ学習」における方策情報の更新周期は、「オンライ

ン学習」よりも長くなる。

[0083] <第3実施形態>

第3実施形態は、より具体的な実施形態に関する。すなわち、第3実施形態は、第2実施形態のバリエーションに関する。

[0084] <制御装置の概要>

図4は、第3実施形態の演算装置80を含む制御装置70の一例を示すブロック図である。図4には、制御装置70の他に、図2と同様に指令実行装置50及び制御対象60が図示されている。

[0085] 制御装置70は、後述するような、「処理フェーズ1」、「処理フェーズ2」、及び「処理フェーズ3」を実行する。制御装置70は、これらの処理を実行することによって、制御対象60の状態 (state) が、より早期に所望状態に近づくよう、方策情報を学習する。

[0086] 方策情報は、制御対象60がある状態である場合に、実行可能なアクションを表す。方策情報は、たとえば、該ある状態と、該アクションとが関連付けられた情報を用いて実現することができる。方策情報は、たとえば、該ある状態を与えた場合に該アクションを算出する処理であってもよい。当該処理は、たとえば、ある関数、または、統計的な手法によって算出された、当該ある状態と、当該アクションとの関係性を表すモデルであってもよい。すなわち、方策情報は、上述した例に限定されない。

[0087] 以降の説明においては、説明の便宜上、制御対象60は、車両や、発電機等（以降、「実環境」と表す）であるとする。あるタイミング（以降、「第1タイミング」と表す）における制御対象60の状態を「第1状態」と表す。あるタイミングの次のタイミング（以降、「第2タイミング」と表す）における制御対象60の状態を「第2状態」と表す。制御対象60の状態は、第1状態に応じたアクションが実施された後に第2状態に変化するとする。また、第1状態と、第2状態とは、必ずしも、相互に異なる状態である必要はなく、同じ状態を表していてもよい。以降の説明においては、説明の便宜上、第1状態と、第2状態との異同によらずに、「第1状態から第2状態に

変化する」と表すこととする。

[0088] 制御装置70は、後述する「処理フェーズ1」にて、複数のタイミングに関して、制御対象60の状態を参照しながら後述するような処理を実行することによって、タイミングごとにアクションを決定する。すなわち、制御装置70は、第1タイミングに関して処理を実行した後に、第2タイミングに関して処理を実行し、さらに、第2タイミングより後のタイミングに関して処理を実行する。したがって、第1タイミング、及び、第2タイミングは、特定のタイミングを表しているわけではなく、制御装置70における処理に関して連続している2つのタイミングを表している。

[0089] (処理フェーズ1)

制御装置70は、第1状態である制御対象60に関して、第1状態と方策情報とに基づきアクションを決定し、決定したアクションを示す制御指令を、指令実行装置50へ出力する。

[0090] 指令実行装置50は、制御装置70によって制御指令を受け取り、受け取った該制御指令が示すアクションを、制御対象60に関して実行する。この結果、制御対象60の状態は、第1状態から第2状態に変化する。

[0091] 説明の便宜上、制御対象60には、制御対象60を観測しているセンサー(図示せず)が取り付けられているとする。センサーは、制御対象60に関して観測した観測値を表すセンサー情報を作成し、作成したセンサー情報を出力するとする。制御対象60を観測しているセンサーは、複数であってもよい。

[0092] 制御装置70は、第1状態に関するアクションが実行された後に、センサーによって作成された該センサー情報を受け取り、受け取った該センサー情報に関する第2状態を推定する。制御装置70は、該第1状態と、該アクションと、該第2状態とが関連付けされた情報(以降、「履歴情報」と表す)を作成する。制御装置70は、作成した履歴情報を、後述する履歴情報記憶部91に格納してもよい。

[0093] 処理フェーズ1に関して上述したような処理が、複数のタイミングに関し

て実行されることにより、後述する履歴情報記憶部41には、複数のタイミングにおける履歴情報が蓄積される。

[0094] (処理フェーズ2)

制御装置70は、処理フェーズ1にて蓄積された履歴情報を用いて、状態遷移情報を更新する(または、作成する)。状態遷移情報を、ニューラルネットワークを用いて作成する場合に、制御装置70は、上述したような履歴情報に含まれているデータを訓練データとして用いて、当該状態遷移情報を作成する。後述するように、制御装置70は、たとえば、構成が相互に異なっているニューラルネットワークを用いて、複数の状態遷移情報を作成する。

[0095] ここで、状態遷移情報について説明する。状態遷移情報は、第1状態と、第2状態との間の関係性を表す情報であり、たとえば、制御対象60の状態遷移(つまり、アクションによる第1状態から第2状態への状態遷移)を、履歴情報を用いてモデル化したものである。すなわち、状態遷移情報を用いることにより、第1状態とアクションとの組み合わせに対応する第2状態を予測することができる。以降、制御対象60の第1状態及び第2状態と区別するために、状態遷移情報の第1状態及び第2状態を、「第1疑似状態」及び「第2疑似状態」と表すことがある。また、「第2疑似状態」を「予測状態」と表すことがある。

[0096] (処理フェーズ3)

制御装置70は、状態遷移情報に基づき、第1疑似状態にて可能な複数の「候補アクション」のそれぞれに関して複数の「予測状態」を決定する。制御装置70は、相互に異なる状態遷移情報(すなわち、各遷移情報ユニット)を用いることによって、複数の第2疑似状態を作成する。

[0097] 状態遷移情報を、ニューラルネットワークを用いて作成する場合に、制御装置70は、第1疑似状態、及び、当該第1疑似状態における候補アクションを表す情報に対して当該状態遷移情報を適用することによって、第2疑似状態を作成する。

- [0098] 処理フェーズ3に関して上述した処理によって、制御装置70は、各候補アクションに関して、複数の予測状態を作成する。制御装置70は、各候補アクションに関して、複数の予測状態のばらつき度を算出する。
- [0099] 制御装置70は、複数の候補アクションの中から、該ばらつき度に基づきアクションを選択する。この選択されたアクションは、後述するように、方策情報の更新に用いられるので、以降、「更新使用アクション」と表すことがある。制御装置70は、複数の候補アクションの中から、算出したばらつき度が上位である候補アクションを特定し、特定した候補アクションの中から更新使用アクションを選択する。制御装置70は、たとえば、複数の候補アクションの中から、算出したばらつき度が最も大きな候補アクションを選択してもよい。
- [0100] 上位は、たとえば、ばらつき度が最も大きいものから、ばらつき度が大きい順に数えて、1%、5%、または、10%等の所定の割合以内であることを表している。
- [0101] 制御装置70は、報酬情報を用いて、1つの候補アクション後の予測状態における報酬程度を求め、求めた報酬程度と、当該1つの候補アクションに対するばらつき度とに基づき、更新使用アクションを選択してもよい。報酬情報は、ある状態が望ましい程度（つまり、「報酬程度」）を表す。報酬情報は、たとえば、該ある状態と、該程度とが関連付けされた情報を用いて実現することができる。報酬情報は、たとえば、該ある状態を与えた場合に該報酬程度を算出する処理であってもよい。当該処理は、たとえば、ある関数、または、統計的な手法によって算出された、当該ある状態と、当該報酬程度との関係性を表すモデルであってもよい。すなわち、報酬情報は、上述した例に限定されない。
- [0102] 予測状態が複数である場合に、制御装置70は、たとえば、各予測状態に関する報酬程度の平均（または、中央値）を求めることによって、候補アクションに関する報酬程度を求める。または、制御装置70は、たとえば、各予測状態の頻度が上位の状態を求め、求めた状態に関する報酬程度の平均（

または、中央値)を求めることによって、候補アクションに関する報酬程度を求める。この場合に、上位は、たとえば、頻度が最も高いものから、頻度が高い順に数えて、1%、5%、または、10%等の所定の割合以内であることを表している。候補アクションに関する報酬程度を求める処理は、上述した例に限定されない。

[0103] また、1つの候補アクションに関する報酬程度と、該1つの候補アクションに関するばらつき度とに基づき、更新使用アクションを選択する処理は、たとえば、該報酬程度と、該ばらつき度とを足し算してもよいし、該報酬程度と、該ばらつき度との重み付き平均を算出してもよい。更新使用アクションを選択する処理は、上述した例に限定されない。

[0104] 制御装置70は、更新使用アクションに基づき、方策情報を更新する。たとえば、制御装置70は、更新使用アクションが処理フェーズ1にて確定的に又は他のアクションに比べて高い確率で選択されるように、方策情報を更新する。この更新された方策情報は、処理フェーズ1にて用いられることになる。

[0105] <制御装置の構成例>

図4において制御装置70は、演算装置80と、記憶装置90とを有している。演算装置30は、状態推定部81と、状態遷移情報更新部(状態遷移情報作成部)82と、制御指令演算部83と、予測状態決定部11と、ばらつき度算出部12と、候補アクション選択部13とを有している。記憶装置90は、履歴情報記憶部91と、状態遷移情報記憶部92と、方策情報記憶部93とを有している。以降、制御装置70の構成を処理フェーズ毎に説明する。

[0106] (処理フェーズ1)

状態推定部81は、制御対象60の状態を表す観測値(パラメタ値、センサー情報)を受け取る。状態推定部81は、受け取った観測値(パラメタ値、センサー情報)に基づき、制御対象60の状態を推定する。

[0107] 制御指令演算部83は、状態推定部81に推定された状態と方策情報記憶

部 9 3 に記憶されている方策情報とに基づきアクションを決定し、決定したアクションを示す制御指令を、指令実行装置 5 0 へ出力する。指令実行装置 5 0 は、制御装置 7 0 によって制御指令を受け取り、受け取った該制御指令が示すアクションを、制御対象 6 0 に関して実行する。この結果、制御対象 6 0 の状態は、第 1 状態から第 2 状態に変化する。

[0108] 状態推定部 8 1 は、制御対象 6 0 の状態（この場合、第 2 状態）を表す観測値（パラメタ値、センサー情報）を受け取る。状態推定部 8 1 は、第 1 状態と、第 1 状態にて実行されたアクションと、該第 2 状態とが関連付けされた履歴情報を作成し、作成した履歴情報を履歴情報記憶部 9 1 に格納する。

[0109] 処理フェーズ 1 に関して上述したような処理を繰り返すことによって、履歴情報記憶部 9 1 には、上記の履歴情報が蓄積される。

[0110] （処理フェーズ 2）

説明の便宜上、ニューラルネットワーク等の統計的な手法（所定の処理手順）を用いて状態遷移情報を作成する例を用いて、処理フェーズ 2 に対応する制御装置 7 0 の構成について説明する。所定の処理手順は、例えば、ニューラルネット等の機械学習法に従った手順である。

[0111] 状態遷移情報更新部 8 2 は、履歴情報記憶部 9 1 に蓄積されている履歴情報を用いて、所定の処理手順に従って、複数の状態遷移情報を作成する。すなわち、状態遷移情報更新部 8 2 は、該履歴情報を訓練データとして、所定の処理手順に従い状態遷移情報を作成し、作成した状態遷移情報を状態遷移情報記憶部 9 2 に格納する。上述したように、状態遷移情報は、第 1 状態と、第 2 状態との関係性を表す。

[0112] たとえば、状態遷移情報更新部 8 2 は、構成が互いに異なる複数のニューラルネットを用いて、複数の遷移情報ユニットを作成してもよい。構成が互いに異なる複数のニューラルネットは、例えば、互いにノードの数又はノード間の接続パターンが異なる複数のニューラルネットである。また、互いに構成の異なる複数のニューラルネットは、あるニューラルネットワークと、当該あるニューラルネットワークにおける一部のノードが存在していない（す

なわち、一部のノードがドロップアウトしている) ニューラルネットワークを用いて実現されていてもよい。

[0113] 状態遷移情報更新部82は、パラメタの初期値が異なる複数のニューラルネットワークを用いて、複数の遷移情報ユニットを作成してもよい。

[0114] 状態遷移情報更新部82は、履歴情報のうちの一部のデータ、または、履歴情報から重複を許してサンプリングしたものを訓練データとして用いてもよい。この場合に、複数の遷移情報ユニットは、相互に異なる訓練データに対して状態遷移情報を作成する。

[0115] なお、所定の処理手順は、ニューラルネットワークに限定されない。例えば所定の処理手順は、SVM (support vector machine)、ランダムフォレスト、バギング (bootstrap aggregating)、又は、ベイジアンネットワークを算出する手順であってもよい。

[0116] (処理フェーズ3)

制御指令演算部83は、第1疑似状態にて可能な複数の候補アクションをそれぞれ示す複数の制御指令を予測状態決定部11へ出力する。

[0117] 予測状態決定部11は、第1疑似状態にて可能な複数の候補アクションと状態遷移情報とに基づき、第1疑似状態にて可能な複数の「候補アクション」のそれぞれに関して複数の予測状態を決定する。制御装置70は、相互に異なる状態遷移情報(すなわち、各遷移情報ユニット)を用いることによって、各候補アクションに関して複数の第2疑似状態を作成する。

[0118] 制御指令演算部83は、予測状態決定部11にて作成された各第2疑似状態を新たな第1疑似状態として、該第1疑似状態にて可能な複数の候補アクションをそれぞれ示す複数の制御指令を予測状態決定部11へ出力する。このとき、制御指令演算部83は、たとえば、予測状態決定部11にて複数の状態遷移情報のうちのある1つを用いて作成された各第2状態情報を新たな第1疑似状態としてもよい。

[0119] 上述したような制御指令演算部83と予測状態決定部11との遣り取りによって、候補アクション選択部13には、第1疑似状態、第2疑似状態、及

び候補アクションの各組み合わせに対応する、ばらつき度が蓄積されることになる。

[0120] ばらつき度算出部 12 は、予測状態決定部 11 によって作成された複数の予測状態のばらつき度（たとえば、分散値、エントロピー等）を算出し、算出したばらつき度を候補アクション選択部 13 へ出力する。ばらつき度は、たとえば、分散値にある数を加えた値等であってもよく、上述した例に限定されない。

[0121] 候補アクション選択部 13 は、複数の候補アクションの中から、該ばらつき度に基づき更新使用アクションを選択する。候補アクション選択部 13 は、たとえば、複数の候補アクションの中から、算出したばらつき度が上位である候補アクションを特定し、特定した候補アクションの中から更新使用アクションを選択する。候補アクション選択部 13 は、たとえば、複数の候補アクションの中から、算出したばらつき度が最も大きな候補アクションを選択してもよい。

[0122] 候補アクション選択部 13 は、更新使用アクションに基づき、方策情報を更新する。たとえば、候補アクション選択部 13 は、処理フェーズ 1 にて制御指令演算部 83 によって更新使用アクションが確定的に又は他の候補アクションに比べて高い確率で選択されるように、方策情報記憶部 93 に記憶されている方策情報を更新する。

[0123] 上述したように候補アクション選択部 13 は、ばらつき度が大きい候補アクションを選択する。ばらつき度は、状態遷移情報に従い算出された結果がばらついていることを表している。このため、ばらつき度が大きい場合には、状態遷移情報が不安定であることを表しているということもできる。すなわち、ばらつき度が大きいアクションを実行することによって、探索が十分でない状態遷移を積極的に探索することができるという効果を奏する。

[0124] 候補アクション選択部 13 は、状態価値情報に基づき、状態に関する価値の程度を表す状態価値情報を作成してもよい。状態価値情報は、たとえば、状態に対して、当該状態の価値の程度を表す関数である。この場合に、価値

は、当該状態を実現することが望ましい程度を表す情報であるとも言えることができる。状態価値情報は、アクション後における制御対象60の状態がどの程度望ましいのかを表す情報ともいうことができる。状態価値情報は、また、当該アクションがどの程度の望ましいのかを表す情報ともいうことができる。

[0125] 候補アクション選択部13は、状態価値情報を作成する処理において、報酬情報を用いてもよい。たとえば、候補アクション選択部13は、各候補アクションに関して算出されたばらつき度を、新たに、状態価値情報として設定してもよい。たとえば、候補アクション選択部13は、各候補アクションに関して算出されたばらつき度を状態価値情報として設定し、その後、当該候補アクションに関する報酬情報を加える等の処理を実行することによって、状態価値情報を更新してもよい。この場合に、ばらつき度は、報酬情報に対する追加的な報酬（疑似追加報酬）であるともいうことができる。

[0126] 状態価値情報を作成する処理は、上述した例に限定されず、たとえば、報酬情報に所定の値を加算した値、報酬情報に所定の値を減算した値、または、報酬情報に所定の値を乗算した値等に基づき実行されてもよい。すなわち、ばらつき度が大きいほど、状態価値情報は、価値が高いことを表す情報であればよい。

[0127] 候補アクション選択部13は、状態価値情報に基づき、複数の候補アクションの中から、該価値の程度が上位の候補アクションを選択し、選択した候補アクションから更新使用アクションを選択してもよい。候補アクション選択部13は、たとえば、算出した価値が最も高い候補アクションを選択してもよい。この場合に、上位は、たとえば、価値の程度が最も高いものから価値の程度が高い順に数えて、1%、5%、または、10%等の所定の割合以内であることを表している。

[0128] <制御装置の動作例>

以上の構成を有する演算装置80の処理動作の一例について説明する。図5は、第3実施形態の演算装置の処理動作の一例を示すフローチャートであ

る。図5に示すフローチャートにおいて、ステップS201は、上記の処理フェーズ1に対応し、ステップS202は、処理フェーズ2に対応し、ステップS203、S204は、処理フェーズ3に対応する。

[0129] 演算装置80は、履歴情報が蓄積されるまで、処理フェーズ1に示された処理を繰り返すことによって、履歴情報を取得する（ステップS201）。

[0130] 演算装置80は、処理フェーズ2に示された処理によって、状態遷移情報を更新する（ステップS202）。

[0131] 演算装置80は、ばらつき度が蓄積されるまで、処理フェーズ3に示された処理によって、ばらつき度を算出する（ステップS203）。

[0132] 演算装置80は、ばらつき度に基づき方策情報を更新する（ステップS204）。そして、処理ステップは、ステップS201（処理フェーズ1）に戻る。

[0133] なお、以上の説明では、演算装置80が、処理フェーズ3にて、ばらつき度を蓄積してから方策情報を更新し、その後直ぐに、処理フェーズ1に戻るものとして説明した。すなわち、以上の説明では、方策情報がバッチ学習されるケースを例にとり説明したが、これに限定されるものではない。例えば、方策情報は、オンライン学習されてもよいし、ミニバッチ学習されてもよい。

[0134] 「オンライン学習」の場合、図5のフローチャートは、ステップS203、S204の処理を繰り返すループとし、該ループが所定回数繰り返されたことを条件に、ステップS201（処理フェーズ1）へ戻るように、修正されてもよい。すなわち、「オンライン学習」の場合、候補アクション選択部13は、ばらつき度が届く度に、方策情報を更新することになる。

[0135] 「ミニバッチ学習」の場合、図5のフローチャートは、「オンライン学習」の場合と同様に、ステップS203、S204の処理を繰り返すループとし、該ループが所定回数繰り返されたことを条件に、ステップS201（処理フェーズ1）へ戻るように、修正されてもよい。ただし、「ミニバッチ学習」の場合、候補アクション選択部13は、「オンライン学習」の場合と異

なり、複数個のばらつき度が蓄積されたタイミングで、方策情報を更新することになる。

[0136] <他の実施形態>

図6は、演算装置のハードウェア構成例を示す図である。図6において演算装置100は、プロセッサ101と、メモリ102とを含んでいる。第1実施形態及び第2実施形態で説明した演算装置10、30、80の状態推定部31、81と、状態遷移情報更新部（状態遷移情報作成部）32、82と、制御指令演算部33、83と、予測状態決定部11と、ばらつき度算出部12と、候補アクション選択部13とは、プロセッサ101がメモリ102に記憶されたプログラムを読み込んで実行することにより実現されてもよい。プログラムは、様々なタイプの非一時的なコンピュータ可読媒体（*non-transitory computer readable medium*）を用いて格納され、演算装置10、30、80に供給することができる。また、プログラムは、様々なタイプの一時的なコンピュータ可読媒体（*transitory computer readable medium*）によって演算装置10、30、80に供給されてもよい。

[0137] 上述したような演算装置は、たとえば、製造工場における各装置を制御する制御装置としても機能することができる。この場合に、各製造工場には、各装置や、製造工場における状態（たとえば、気温、湿度、視界等）等を測定するセンサーが配置される。各センサーは、各装置や、製造工場における状態等を測定し、測定した状態を表す観測情報を作成する。この場合に、観測情報は、製造工場において観測される状態を表す情報である。

[0138] 演算装置は、当該観測情報を受け取り、上述したような処理を行うことにより決定されたアクションに従い、各装置を制御する。たとえば、装置が、材料の量を調整するバルブである場合に、演算装置は、決定したアクションに従い、バルブを閉める、または、バルブを開ける等の制御を行う。または、装置が、温度を調整するヒータである場合に、演算装置は、決定したアクションに従い、設定温度を上げる、または、設定温度を下げる等の制御を行

う。

[0139] 製造工場における各装置を制御する例を参照しながら、制御例について説明したが、制御例は、上述した例に限定されない。たとえば、演算装置は、上述したような処理と同様な処理を行うことによって、化学工場における各装置を制御する制御装置、発電所における各装置を制御する制御装置としても機能することができる。

[0140] 以上、実施の形態を参照して本願発明を説明したが、本願発明は上記によって限定されるものではない。本願発明の構成や詳細には、発明のScope内で当業者が理解し得る様々な変更をすることができる。

符号の説明

- [0141] 10, 30, 80 演算装置 (アクション決定装置)
- 11 予測状態決定部
 - 12 ばらつき度算出部
 - 13 候補アクション選択部
 - 20, 70 制御装置
 - 31, 81 状態推定部
 - 32, 82 状態遷移情報更新部 (状態遷移情報作成部)
 - 33, 83 制御指令演算部
 - 40, 90 記憶装置
 - 41, 91 履歴情報記憶部
 - 42, 92 状態遷移情報記憶部
 - 43, 93 方策情報記憶部
 - 50 指令実行装置
 - 60 制御対象

請求の範囲

- [請求項1] 第1 タイミングでの第1 状態と、前記第1 タイミング以降の第2 タイミングでの第2 状態との関係性を表す遷移情報を複数用いて、第1 状態にて可能な複数の候補アクションのそれぞれに関して複数の第2 状態を決定する決定手段と、
- 各前記候補アクションに関して、前記複数の第2 状態のばらつき度を算出する算出手段と、
- 前記ばらつき度に基づき、前記複数の候補アクションから一部の候補アクションを選択する選択手段と、
- を具備する演算装置。
- [請求項2] 前記選択手段は、前記複数の候補アクションの中から、前記一部の候補アクションとして、前記ばらつき度が上位の前記候補アクションを選択する、
- 請求項1 に記載の演算装置。
- [請求項3] 前記選択手段は、前記一部の候補アクションとして、前記ばらつき度が最大の前記候補アクションを選択する、
- 請求項1 に記載の演算装置。
- [請求項4] 2つの状態と、該2つの状態間におけるアクションとが関連付けされたセットを含む履歴情報に基づき、所定の処理手順に従い、前記遷移情報を作成する作成手段をさらに具備する、
- 請求項1 乃至請求項3 のいずれかに記載の演算装置。
- [請求項5] 前記所定の処理手順は、ニューラルネットを算出する手順である、
- 請求項4 に記載の演算装置。
- [請求項6] 前記作成手段は、複数の前記遷移情報を、互いに構成の異なる複数の前記ニューラルネットを用いて作成する、
- 請求項5 に記載の演算装置。
- [請求項7] 前記作成手段は、複数の前記遷移情報を、パラメタの初期値が異なる複数の前記ニューラルネットを用いて作成する、

請求項5に記載の演算装置。

[請求項8] 複数の前記遷移情報を、前記履歴情報のうちの互いに異なるセットを複数の前記ニューラルネットに入力することによって作成する、

請求項5に記載の演算装置。

[請求項9] 情報処理装置によって、第1タイミングでの第1状態と、前記第1タイミング以降の第2タイミングでの第2状態との関係性を表す遷移情報を複数用いて、第1状態にて可能な複数の候補アクションのそれぞれに関して複数の第2状態を決定し、

各前記候補アクションに関して、前記複数の第2状態のばらつき度を算出し、

前記ばらつき度に基づき、前記複数の候補アクションから一部の候補アクションを選択する、

アクション決定方法。

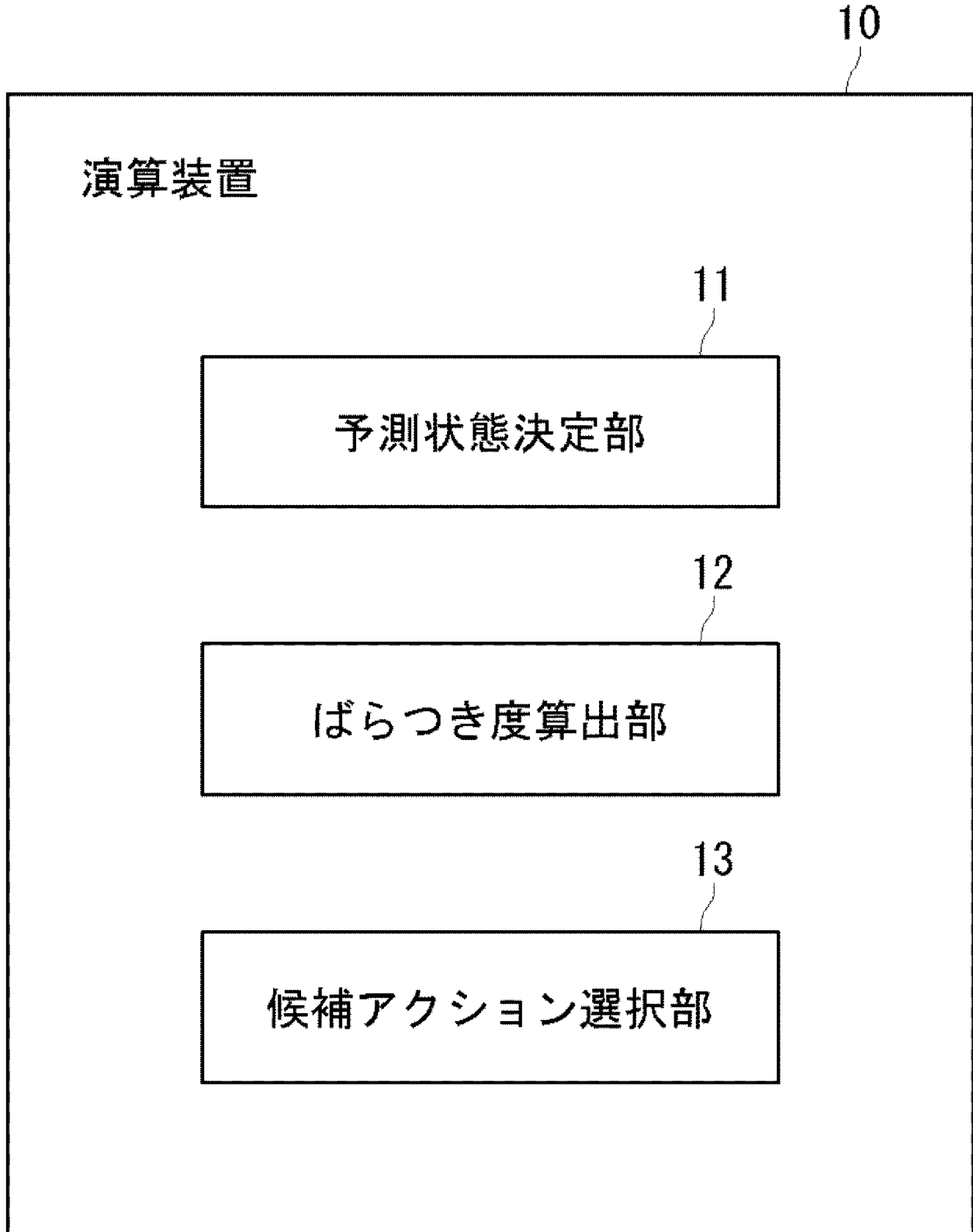
[請求項10] 第1タイミングでの第1状態と、前記第1タイミング以降の第2タイミングでの第2状態との関係性を表す遷移情報を複数用いて、第1状態にて可能な複数の候補アクションのそれぞれに関して複数の第2状態を決定し、

各前記候補アクションに関して、前記複数の第2状態のばらつき度を算出し、

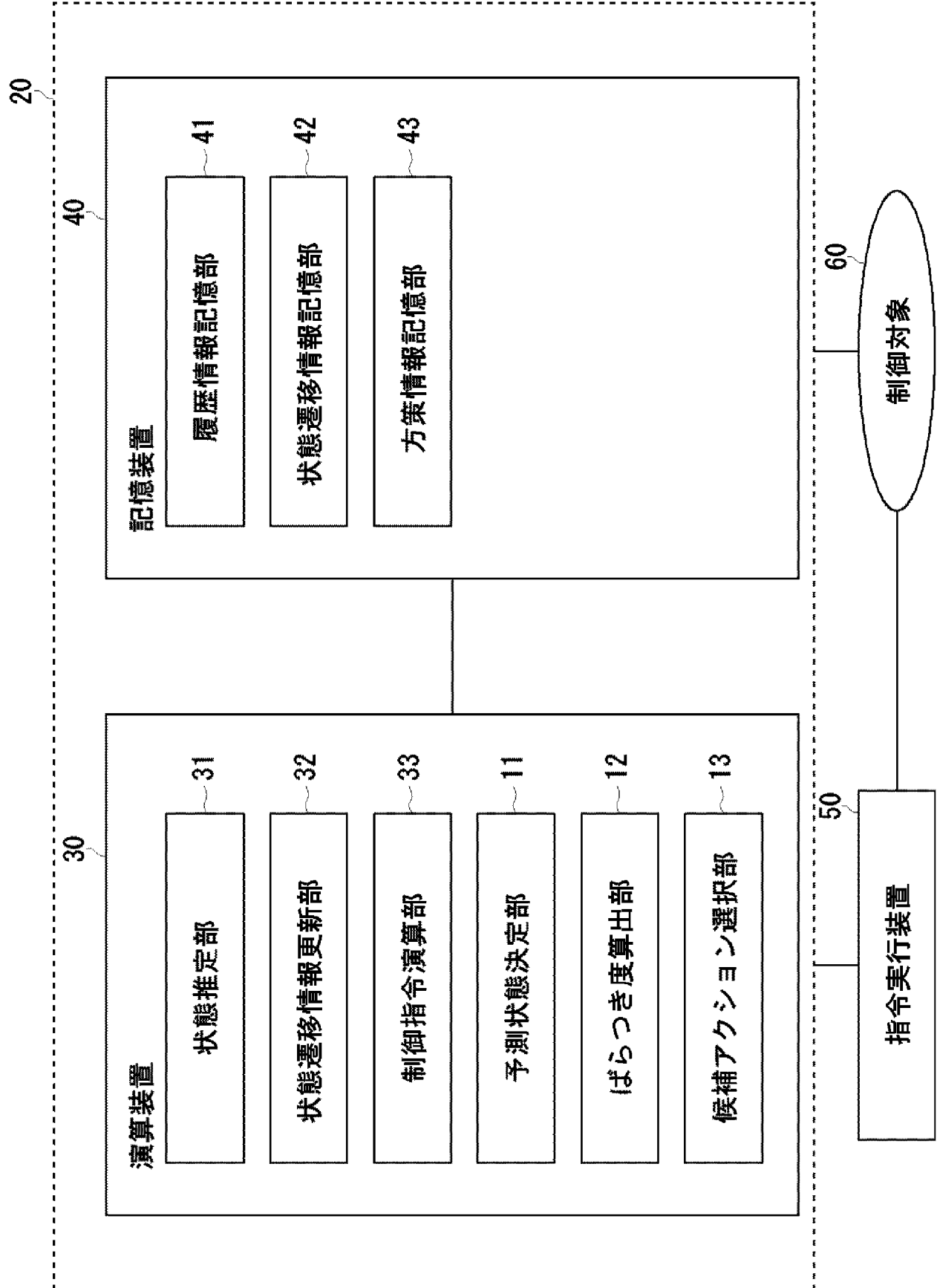
前記ばらつき度に基づき、前記複数の候補アクションから一部の候補アクションを選択する、

処理を、演算装置に実行させる制御プログラムを格納する非一時的なコンピュータ可読媒体。

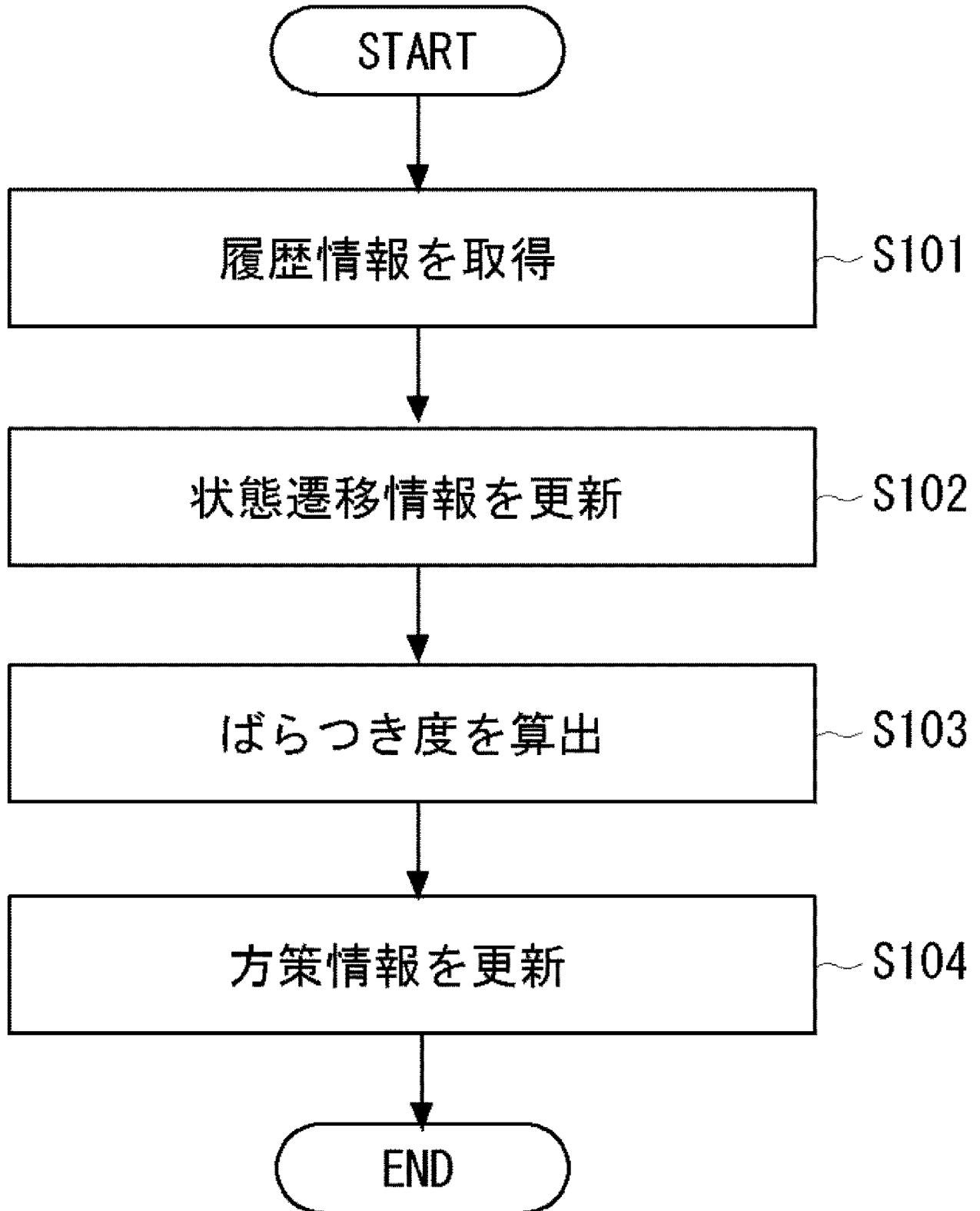
[図1]



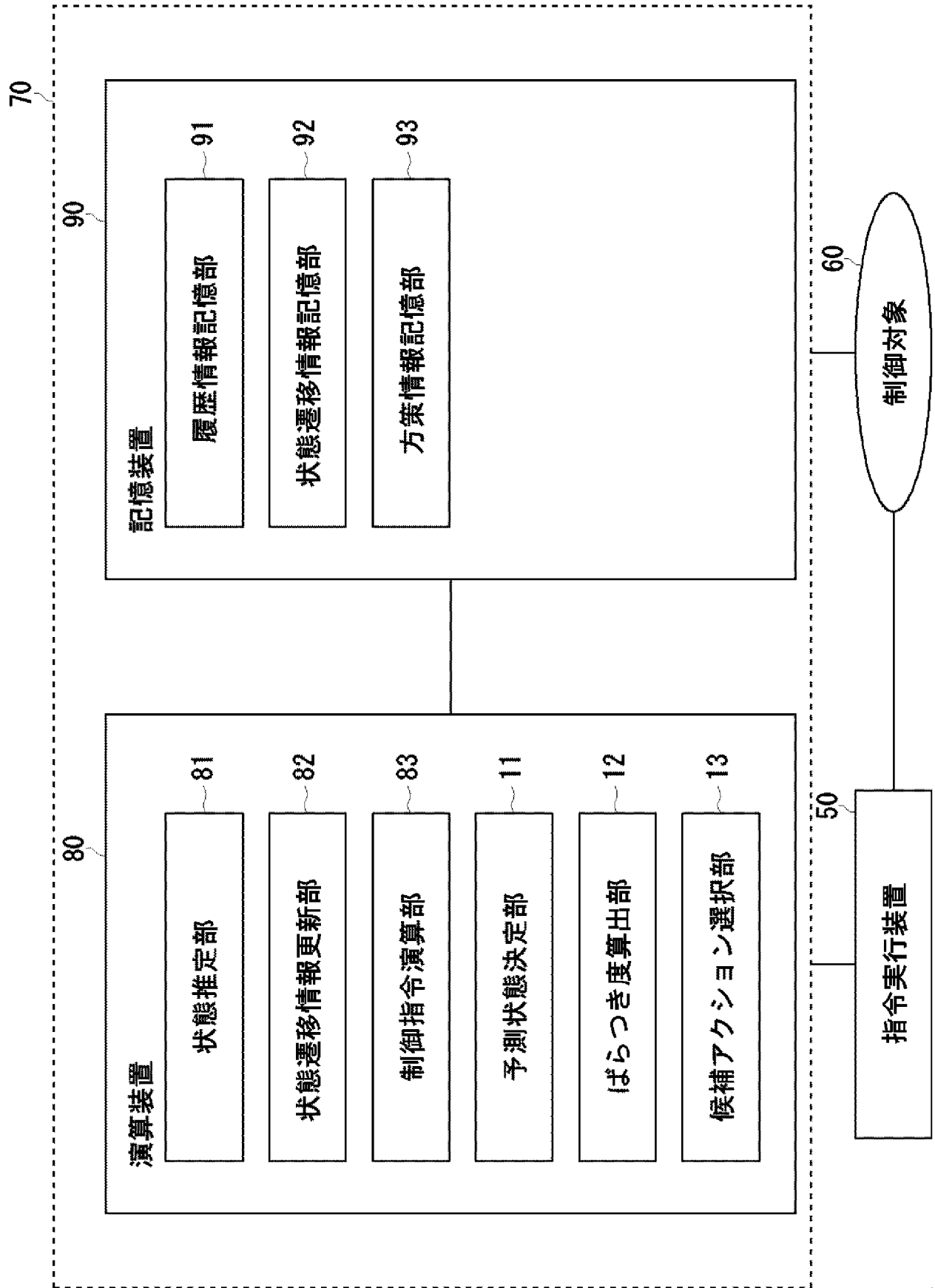
[図2]



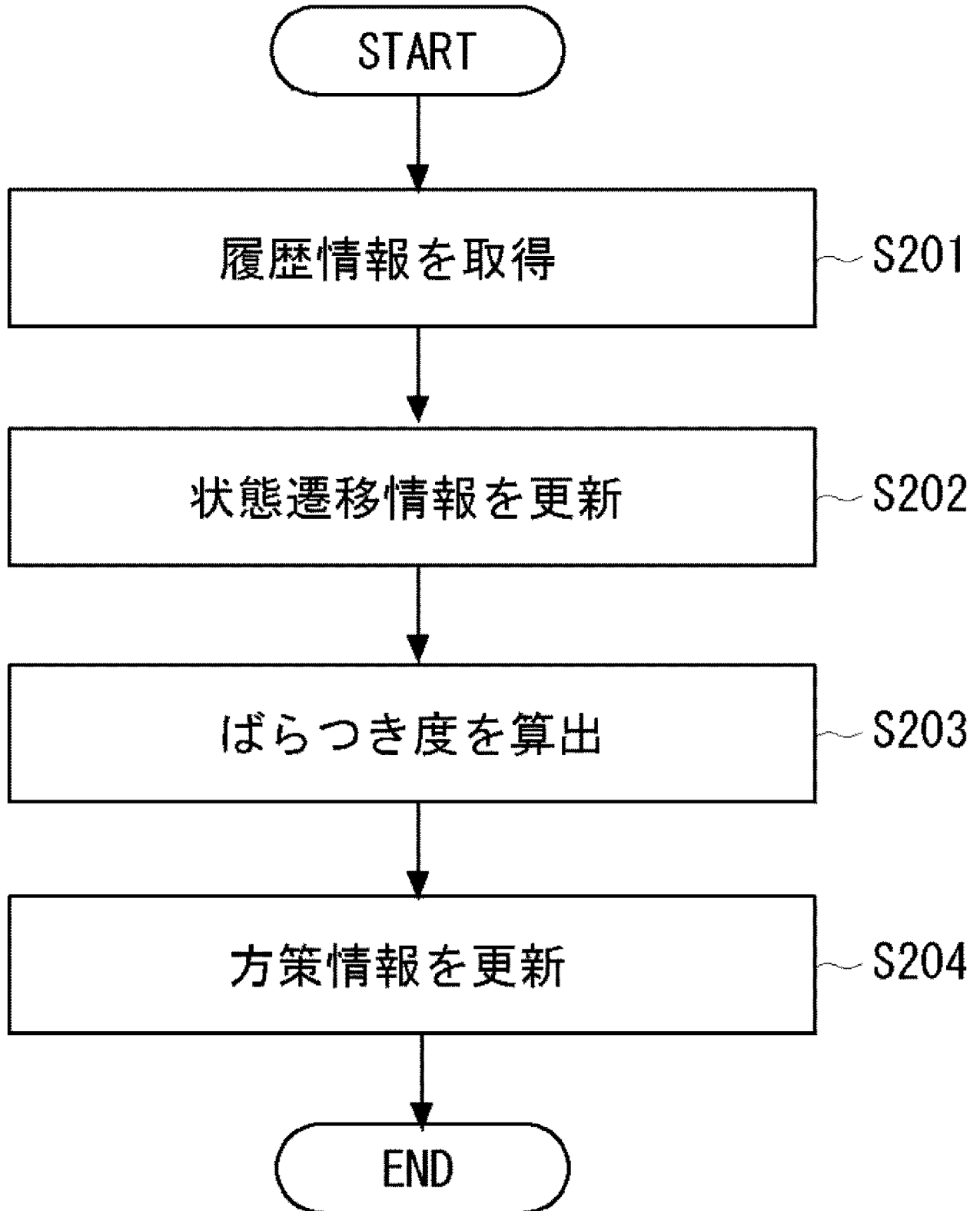
[図3]



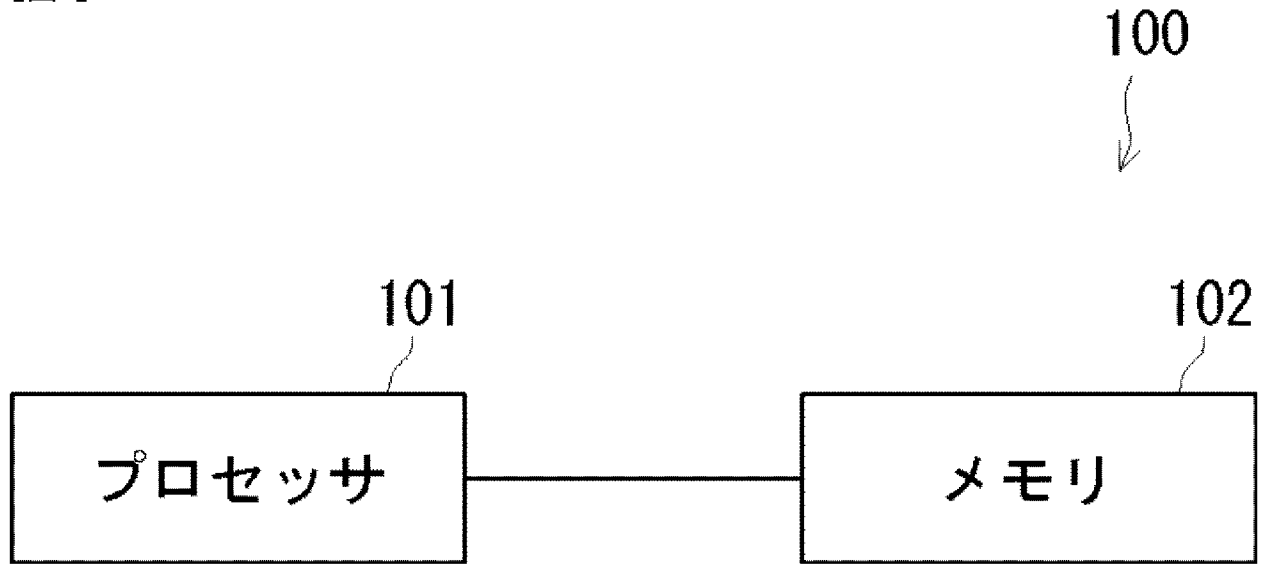
[図4]



[図5]



[図6]



INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP2018/045947

A. CLASSIFICATION OF SUBJECT MATTER Int. Cl. G06N20/00 (2019.01) i According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED Minimum documentation searched (classification system followed by classification symbols) Int. Cl. G06N3/00-99/00 Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched Published examined utility model applications of Japan 1922-1996 Published unexamined utility model applications of Japan 1971-2019 Registered utility model specifications of Japan 1996-2019 Published registered utility model applications of Japan 1994-2019 Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y A	齋藤雅矩ほか, エージェントの行動履歴の活用による Q-learning の学習効率向上, 電気学会研究会資料, 07 December 2014, pp. 29-34, in particular, p. 31, right column, line 9 to p. 32, left column, line 41 (SAITO, Masanori et al. Improving efficiency of Q-learning by using the agent's action history. IEEJ Working Group Materials.)	1-4, 9-10 5-8
Y A	HAARNOJA, Tuomas et al., Reinforcement Learning with Deep Energy-Based Policies, arXiv [online], 21 July 2017, [retrieved on 28 January 2019], Retrieved from the Internet: <URL: https://arxiv.org/pdf/1702.08165v2 >, in particular, page 2, right column, line 19 to page 3, left column, line 27	1-4, 9-10 5-8
<input type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
* Special categories of cited documents: "A" document defining the general state of the art which is not considered to be of particular relevance "E" earlier application or patent but published on or after the international filing date "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) "O" document referring to an oral disclosure, use, exhibition or other means "P" document published prior to the international filing date but later than the priority date claimed "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art "&" document member of the same patent family		
Date of the actual completion of the international search 28.01.2019		Date of mailing of the international search report 12.02.2019
Name and mailing address of the ISA/ Japan Patent Office 3-4-3, Kasumigaseki, Chiyoda-ku, Tokyo 100-8915, Japan		Authorized officer Telephone No.

A. 発明の属する分野の分類（国際特許分類（IPC）） Int.Cl. G06N20/00(2019.01) i										
B. 調査を行った分野 調査を行った最小限資料（国際特許分類（IPC）） Int.Cl. G06N3/00-99/00										
最小限資料以外の資料で調査を行った分野に含まれるもの <table border="0"> <tr> <td>日本国実用新案公報</td> <td>1922-1996年</td> </tr> <tr> <td>日本国公開実用新案公報</td> <td>1971-2019年</td> </tr> <tr> <td>日本国実用新案登録公報</td> <td>1996-2019年</td> </tr> <tr> <td>日本国登録実用新案公報</td> <td>1994-2019年</td> </tr> </table>			日本国実用新案公報	1922-1996年	日本国公開実用新案公報	1971-2019年	日本国実用新案登録公報	1996-2019年	日本国登録実用新案公報	1994-2019年
日本国実用新案公報	1922-1996年									
日本国公開実用新案公報	1971-2019年									
日本国実用新案登録公報	1996-2019年									
日本国登録実用新案公報	1994-2019年									
国際調査で使用した電子データベース（データベースの名称、調査に使用した用語）										
C. 関連すると認められる文献										
引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求項の番号								
Y A	齋藤雅矩ほか，エージェントの行動履歴の活用によるQ-learningの学習効率向上，電気学会研究会資料，2014.12.07，pp.29-34，特に第31頁右欄第9行～第32頁左欄第41行	1-4，9-10 5-8								
Y A	HAARNOJA, Tuomas et al., Reinforcement Learning with Deep Energy-Based Policies, arXiv [online], 2017.07.21, [retrieved on 2019.01.28], Retrieved from the Internet: <URL: https://arxiv.org/pdf/1702.08165v2>, 特に第2頁右欄第19行～第3頁左欄第27行	1-4，9-10 5-8								
<input type="checkbox"/> C欄の続きにも文献が列挙されている。 <input type="checkbox"/> パテントファミリーに関する別紙を参照。										
* 引用文献のカテゴリー 「A」特に関連のある文献ではなく、一般的技術水準を示すもの 「E」国際出願日前の出願または特許であるが、国際出願日以後に公表されたもの 「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献（理由を付す） 「O」口頭による開示、使用、展示等に言及する文献 「P」国際出願日前で、かつ優先権の主張の基礎となる出願										
の日の後に公表された文献 「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの 「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの 「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの 「&」同一パテントファミリー文献										
国際調査を完了した日 28.01.2019	国際調査報告の発送日 12.02.2019									
国際調査機関の名称及びあて先 日本国特許庁（ISA/J P） 郵便番号100-8915 東京都千代田区霞が関三丁目4番3号	特許庁審査官（権限のある職員） 北元 健太 電話番号 03-3581-1101 内線 3545	5B 3856								