

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号
特許第7623879号
(P7623879)

(45)発行日 令和7年1月29日(2025.1.29)

(24)登録日 令和7年1月21日(2025.1.21)

(51)国際特許分類 F I
G 0 6 N 20/00 (2019.01) G 0 6 N 20/00

請求項の数 6 (全40頁)

(21)出願番号	特願2021-67848(P2021-67848)	(73)特許権者	000005108 株式会社日立製作所 東京都千代田区丸の内一丁目6番6号
(22)出願日	令和3年4月13日(2021.4.13)	(74)代理人	110001678 藤央弁理士法人
(65)公開番号	特開2022-162824(P2022-162824 A)	(72)発明者	森澤 利浩 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
(43)公開日	令和4年10月25日(2022.10.25)	(72)発明者	古林 雅佳 東京都千代田区丸の内一丁目6番6号 株式会社日立製作所内
審査請求日	令和6年2月22日(2024.2.22)	審査官	千葉 久博

最終頁に続く

(54)【発明の名称】 行動制御計画装置、行動制御計画方法及び行動制御計画システム

(57)【特許請求の範囲】

【請求項1】

プロセッサとメモリを有して、他移動体の行動に対して自移動体の行動を計画する行動制御計画装置であって、

前記自移動体と他移動体の行動をシミュレーション条件として設定するシミュレーション条件設定部と、

前記シミュレーション条件と予め設定された機械学習モデルに基づいて、前記自移動体と他移動体の位置関係を所定の時間間隔毎にステップデータとして出力するシミュレータと、

前記ステップデータを取得して、前記機械学習モデルに前記ステップデータを与えて学習結果として前記自移動体と他移動体の前記位置関係を示す状態量と報酬と状態を含むエピソードデータを出力させ、当該エピソードデータをエピソードデータ蓄積部に蓄積する強化学習部と、

前記エピソードデータ蓄積部に蓄積された前記エピソードデータを分析して学習パターンを生成するエピソード分析部と、

前記学習パターンに該当する前記エピソードデータを前記エピソードデータ蓄積部から取得して学習用エピソードデータを生成し、前記機械学習モデルに前記学習用エピソードデータを与えて学習させるエピソードデータ学習部と、
を有し、

前記強化学習部は、

10

20

前記ステップデータの開始から所定の終了条件に達するまでのステップデータを1つのエピソードとし、前記エピソードの前記ステップデータに対応するステップを前記エピソードデータに設定し、前記エピソードデータの前記状態には前記ステップデータ毎のシミュレーション結果を設定し、前記エピソードには前記シミュレーション結果に応じて所定の目的を達成したか否かを示すエピソード種別を設定し、

前記エピソード分析部は、

前記エピソード種別と前記シミュレーション結果と前記エピソードの数と前記エピソードデータのステップ数に基づいて前記学習パターンを生成し、

前記エピソード分析部は、

前記シミュレーション結果毎に前記エピソードの数の範囲と、前記エピソードデータのステップ数の範囲を設定し、

10

前記エピソードデータ学習部は、

前記エピソードデータ蓄積部から、前記シミュレーション結果毎に前記エピソードの数の範囲と、前記エピソードデータのステップ数の範囲に該当するエピソードを取得して前記機械学習モデルに学習させ、

前記エピソードデータ学習部は、

前記シミュレーション結果毎に、前記エピソードの数とステップ数の範囲で抽出したエピソードを、前記エピソードの順序で所定の置換を行う、ことを特徴とする行動制御計画装置。

【請求項2】

20

請求項1に記載の行動制御計画装置であって、

前記エピソードデータ学習部は、

前記取得したエピソードに対応する前記エピソードデータを前記機械学習モデルに直接入力して学習させることを特徴とする行動制御計画装置。

【請求項3】

請求項1に記載の行動制御計画装置であって、

前記所定の目的は、

前記自移動体が前記他移動体の追跡であることを特徴とする行動制御計画装置。

【請求項4】

請求項1に記載の行動制御計画装置であって、

30

前記エピソード分析部は、

前記エピソードの数の範囲と、前記エピソードデータのステップ数の範囲において、前記シミュレーション結果の種別に応じた所定の比率で前記エピソードを割り当てることを特徴とする行動制御計画装置。

【請求項5】

プロセッサとメモリを有する計算機が、他移動体の行動に対して自移動体の行動を計画する行動制御計画方法であって、

前記計算機が、前記自移動体と他移動体の行動をシミュレーション条件として設定するシミュレーション条件設定ステップと、

前記計算機が、前記シミュレーション条件と予め設定された機械学習モデルに基づいて、前記自移動体と他移動体の位置関係を所定の時間間隔毎にステップデータとして出力するシミュレーションステップと、

40

前記計算機が、前記ステップデータを取得して、前記機械学習モデルに前記ステップデータを与えて学習結果として前記自移動体と他移動体の前記位置関係を示す状態量と報酬と状態を含むエピソードデータを出力させ、当該エピソードデータをエピソードデータ蓄積部に蓄積する強化学習ステップと、

前記計算機が、前記エピソードデータ蓄積部に蓄積された前記エピソードデータを分析して学習パターンを生成するエピソード分析ステップと、

前記計算機が、前記学習パターンに該当する前記エピソードデータを前記エピソードデータ蓄積部から取得して学習用エピソードデータを生成し、前記機械学習モデルに前記学

50

習用エピソードデータを与えて学習させるエピソードデータ学習ステップと、
を含み、

前記強化学習ステップは、

前記ステップデータの開始から所定の終了条件に達するまでのステップデータを1つのエピソードとし、前記エピソードの前記ステップデータに対応するステップを前記エピソードデータに設定し、前記エピソードデータの前記状態には前記ステップデータ毎のシミュレーション結果を設定し、前記エピソードには前記シミュレーション結果に応じて所定の目的を達成したか否かを示すエピソード種別を設定し、

前記エピソード分析ステップは、

前記エピソード種別と前記シミュレーション結果と前記エピソードの数と前記エピソードデータのステップ数に基づいて前記学習パターンを生成し

10

前記エピソード分析ステップは、

前記シミュレーション結果毎に前記エピソードの数の範囲と、前記エピソードデータのステップ数の範囲を設定し、

前記エピソードデータ学習ステップは、

前記エピソードデータ蓄積部から、前記シミュレーション結果毎に前記エピソードの数の範囲と、前記エピソードデータのステップ数の範囲に該当するエピソードを取得して前記機械学習モデルに学習させ、

前記エピソードデータ学習ステップは、

前記シミュレーション結果毎に、前記エピソードの数とステップ数の範囲で抽出したエピソードを、前記エピソードの順序で所定の置換を行う、ことを特徴とする行動制御計画方法。

20

【請求項6】

プロセッサとメモリを有する行動制御計画装置と、

前記行動制御計画装置に接続されて所定のシミュレーションを行うシミュレータと、を有して、他移動体の行動に対して自移動体の行動を計画する行動制御計画システムであって、

前記シミュレータは、

前記行動制御計画装置から受け付けたシミュレーション条件で前記シミュレーションを実施し、

30

前記行動制御計画装置は、

前記自移動体と他移動体の行動をシミュレーション条件として設定し、前記シミュレーション条件と予め設定された機械学習モデルを前記シミュレータに送信して、前記自移動体と他移動体の位置関係を所定の時間間隔毎にステップデータとして出力させるシミュレーション条件設定部と、

前記ステップデータを取得して、前記機械学習モデルに前記ステップデータを与えて学習結果として前記自移動体と他移動体の前記位置関係を示す状態量と報酬と状態を含むエピソードデータを出力させ、当該エピソードデータをエピソードデータ蓄積部に蓄積する強化学習部と、

前記エピソードデータ蓄積部に蓄積された前記エピソードデータを分析して学習パターンを生成するエピソード分析部と、

40

前記学習パターンに該当する前記エピソードデータを前記エピソードデータ蓄積部から取得して学習用エピソードデータを生成し、前記機械学習モデルに前記学習用エピソードデータを与えて学習させるエピソードデータ学習部と、

を有し、

前記強化学習部は、

前記ステップデータの開始から所定の終了条件に達するまでのステップデータを1つのエピソードとし、前記エピソードの前記ステップデータに対応するステップを前記エピソードデータに設定し、前記エピソードデータの前記状態には前記ステップデータ毎のシミュレーション結果を設定し、前記エピソードには前記シミュレーション結果に応じて所定

50

の目的を達成したか否かを示すエピソード種別を設定し、
 前記エピソード分析部は、
 前記エピソード種別と前記シミュレーション結果と前記エピソードの数と前記エピソードデータのステップ数に基づいて前記学習パターンを生成し、
 前記エピソード分析部は、
 前記シミュレーション結果毎に前記エピソードの数の範囲と、前記エピソードデータのステップ数の範囲を設定し、
 前記エピソードデータ学習部は、
 前記エピソードデータ蓄積部から、前記シミュレーション結果毎に前記エピソードの数の範囲と、前記エピソードデータのステップ数の範囲に該当するエピソードを取得して前記機械学習モデルに学習させ、
 前記エピソードデータ学習部は、
 前記シミュレーション結果毎に、前記エピソードの数とステップ数の範囲で抽出したエピソードを、前記エピソードの順序で所定の置換を行う、ことを特徴とする行動制御計画システム。

10

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、艦艇、車両といった移動体の行動において、特に他の移動体との位置関係を鑑みて行動する追尾や哨戒といった行動において、他移動体に対する自移動体の行動制御方法を自動的に計画する行動制御計画装置に関する。

20

【背景技術】

【0002】

環境に応じて移動体の移動方法、すなわち行動制御内容を決めるためには、様々な目的、要件を鑑みることが必要である。環境としては、移動範囲の制限、障害物といった固定的な制限の他にも、他の移動体との位置関係による制限もある。特に、他の移動体との位置関係については、それが障害物であって回避すればよいといった目的の他にも、他の移動体との距離を一定に保つ、もしくは追跡、追尾する、また他の移動体を探索、検知するといった目的もある。

【0003】

30

制御対象である自移動体 (own vehicle) の行動制御内容を決めるためには、従来は制御理論などにに基づき制御方法が設計されるが、強化学習を利用することで行動制御内容を計画できる。近年、状態量と行動 (方策)、行動価値もしくは報酬の関係で定義される、深層ニューラルネットワークを訓練することで強化学習を行う、深層強化学習といった方法により学習能力が向上しており、強化学習技術の適用が検討される。

【0004】

特許文献1では、障害物が設定された環境において移動物の行動を決定するために、深層強化学習により、移動物の行動を決定する技術が提案される。

【0005】

特許文献2では、特に深層強化学習技術に関して、深層ニューラルネットワークを訓練する際、ステップ毎の状態量、行動、報酬のデータに対し、学習の効果が高くなる指標である期待学習進展度を計算して関連付け、リプレイメモリ (訓練のためのステップデータの蓄積) から優先順位付けして訓練に用いる技術が提案される。

40

【先行技術文献】

【特許文献】

【0006】

【文献】特開2018-198012号公報

【文献】特開2020-47285号公報

【発明の概要】

【発明が解決しようとする課題】

50

【0007】

環境に応じて移動体の行動制御内容を強化学習技術を利用して決定する場合、学習を実現することが問題となる。移動体の行動の目的を達成するための前提とする環境などの条件が単純であれば、学習は容易であるが、条件が複雑であれば、学習にかかる訓練の回数が多く必要となり、さらに学習自体を達成できなくなる。

【0008】

特に移動体の行動制御計画では、追跡、探索といった目的もあり、他の移動体（以下、他移動体）との位置関係を反映することが必要となる。この場合、自移動体の位置、移動方向に対して、近づく、離れる、横切るなど、他移動体の位置、移動方向が多様であるが、いずれに対しても適切な行動をとれることが重要である。

10

【0009】

特許文献1では、環境は動かない障害物、もしくは通り道といった設定であり、移動体の一つの学習結果から他の多様な環境に対応する学習方法が示されない。

【0010】

特許文献2では、深層ニューラルネットワークの訓練を効率化して、学習を成功させるための技術である。環境に応じた学習をするためには、成功が失敗に至った時点までを示すエピソード自体の設定を変えて学習しなければならないが、その設定方法に関する工夫は示されない。

【0011】

本発明の目的は、多様な環境に対して、特に、自移動体に対して他移動体が様々な行動を行うことに対応して、自移動体が行動できるように行動を計画することである。この目的のため、他移動体の異なる行動に対して強化学習を行った結果のエピソードのデータを利用して、エピソードデータを用いて自移動体の行動計画を学習する技術を提供するものである。

20

【課題を解決するための手段】

【0012】

本発明は、プロセッサとメモリを有して、他移動体の行動に対して自移動体の行動を計画する行動制御計画装置であって、前記自移動体と他移動体の行動をシミュレーション条件として設定するシミュレーション条件設定部と、前記シミュレーション条件と予め設定された機械学習モデルに基づいて、前記自移動体と他移動体の位置関係を所定の時間間隔毎にステップデータとして出力するシミュレータと、前記ステップデータを取得して、前記機械学習モデルに前記ステップデータを与えて学習結果として前記自移動体と他移動体の前記位置関係を示す状態量と報酬と状態を含むエピソードデータを出力させ、当該エピソードデータをエピソードデータ蓄積部に蓄積する強化学習部と、前記エピソードデータ蓄積部に蓄積された前記エピソードデータを分析して学習パターンを生成するエピソード分析部と、前記学習パターンに該当する前記エピソードデータを前記エピソードデータ蓄積部から取得して学習用エピソードデータを生成し、前記機械学習モデルに前記学習用エピソードデータを与えて学習させるエピソードデータ学習部と、を有する。

30

【発明の効果】

【0013】

本発明によると、様々な環境に対して、特に自移動体に対して他移動体の行動が多様な関係となる場合でも、例えば追跡、検知といった目的を達成するための、自移動体の行動を学習できる。

40

【0014】

本明細書において開示される主題の、少なくとも一つの実施の詳細は、添付されている図面と以下の記述の中で述べられる。開示される主題のその他の特徴、態様、効果は、以下の開示、図面、請求項により明らかにされる。

【図面の簡単な説明】

【0015】

【図1】本発明の実施例を示し、行動制御計画装置で行われる処理の一例を示した図であ

50

る。

【図 2】本発明の実施例を示し、強化学習と、エピソード分析及びエピソードデータ学習による行動制御計画の処理の一例を示した図である。

【図 3】本発明の実施例を示し、移動体の行動制御計画のフローチャートの一例を示した図である。

【図 4】本発明の実施例を示し、移動体の行動制御計画を実行する行動制御計画装置の構成の一例を示した図である。

【図 5】本発明の実施例を示し、追跡における自移動体、他移動体の行動環境の例を示した図である。

【図 6】本発明の実施例を示し、哨戒における自移動体、他移動体の行動環境の例を示した図である。

10

【図 7 A】本発明の実施例を示し、哨戒領域の構成の例を示した図である。

【図 7 B】本発明の実施例を示し、他移動体の行動パターンの例を示した図である。

【図 8】本発明の実施例を示し、方策関数、状態価値関数の深層ニューラルネットワークによる構成の例を示した図である。

【図 9 A】本発明の実施例を示し、エピソードデータの構成の失敗例を示した図である。

【図 9 B】本発明の実施例を示し、エピソードデータの構成の成功例を示した図である。

【図 10】本発明の実施例を示し、一連の強化学習の手順の、フローチャートの一例を示した図である。

【図 11A】本発明の実施例を示し、1つの強化学習過程についてのエピソードデータの例を示した図である。

20

【図 11B】本発明の実施例を示し、1つの強化学習過程についてのエピソードデータの例を示した図である。

【図 12】本発明の実施例を示し、多くの成功した強化学習の学習過程についての、ステップ数のプロットの例を示した図である。

【図 13】本発明の実施例を示し、強化学習過程についての、失敗、成功エピソード別のステップ数の関係の例を示した図である。

【図 14】本発明の実施例を示し強化学習過程における失敗、成功エピソードの出現の範囲の一例を示した図である。

【図 15】本発明の実施例を示し、学習パターンの一例を示した図である。

30

【図 16】本発明の実施例を示し、強化学習過程に対するステップ数のプロットとして、学習パターンのグラフの一例を示した図である。

【図 17】本発明の実施例を示し、学習パターンに対応するエピソードデータのサンプリング結果の一例を示した図である。

【図 18 A】本発明の実施例を示し、エピソードデータのサンプリング結果を置換した例を示した図である。

【図 18 B】本発明の実施例を示し、エピソードデータのサンプリング結果を置換した他の例を示した図である。

【図 19】本発明の実施例を示し、一連のエピソードデータ学習の処理手順の、フローチャートの一例を示した図である。

40

【発明を実施するための形態】

【0016】

以下、本発明の実施形態を添付図面に基づいて説明する。

【0017】

自律的に移動する移動体（自移動体）の行動制御計画のためには、自移動体と他移動体を移動させるシミュレーションによる強化学習技術を利用する。強化学習では、自移動体と他移動体の位置関係の初期値から、特定の目的の達成について成功、もしくは失敗するまで1ステップずつ行動する。なお、ステップは所定の時間間隔毎に実行するシミュレーション（又は学習や認識）の処理を示す。

【0018】

50

自移動体と他移動体を移動させるシミュレーションが開始されてから、自移動体の行動が成功又は失敗に至った時点までを1つのエピソードと言う。行動制御計画装置は、エピソードを繰り返すことで目的を達成する行動制御を機械学習モデルで学習する。強化学習のためには状態量と行動（方策）、状態価値の関係は深層ニューラルネットワーク（機械学習モデル）で定義され、行動制御計画装置では、エピソード毎に状態、行動、報酬のステップ毎のデータを用いて深層ニューラルネットワークが強化学習とエピソードデータ学習によって訓練される。

【0019】

なお、エピソードを構成するステップは、例えば、所定時間間隔毎の状態量と、行動と、状態価値が含まれる。そして、状態量には、自移動体や他移動体の位置や速度を含むことができ、行動には自移動体の進路や速度を含むことができ、状態価値には、ステップ毎の報酬を含むことができる。なお、他移動体は1つ以上の移動体で、複数の移動体も対象とすることができる。

10

【0020】

状態量と行動（方策）、状態価値との関係は損失関数を最小化する最適化により訓練されるので、最適化を効率化するステップデータを優先付けて、選別する。またステップの時系列の並びであるエピソードは、学習完了までの順序の並びとなる。なお、損失関数は、深層ニューラルネットワークの予測精度を評価する関数であり、例えば、回帰モデル等の周知又は公知の技術を適用することができる。

【0021】

自移動体に対する他移動体の行動パターンが単純、もしくは一通りである場合、強化学習を適用すれば学習は達成される。しかし他移動体の行動パターンは様々であることが想定され、一通りであることを仮定できない。

20

【0022】

そこで行動制御計画装置は、他移動体と自移動体の多様な行動パターンについて、各行動パターン別に強化学習を行ったときのエピソードデータを蓄積する。エピソードデータとは、ステップ毎の状態量、行動（方策）、状態価値の並びで構成されたデータである。また、エピソードデータは、複数のステップデータを含むことができる。本実施例では予め蓄積したエピソードデータより、学習が成功に至るまでのエピソード毎の学習過程を分析する。

30

【0023】

行動制御計画装置は、分析の結果として、ステップ数に応じた失敗のエピソードと成功のエピソードの、学習開始からのエピソードの順番である、学習パターンを設定する。なお、行動制御計画装置は、1つのエピソードに含まれるステップ数の多寡に応じて学習パターンを設定することができる。

【0024】

状態量を入力とし、行動（方策）、状態価値（報酬）を出力として、それらの関係が定義される深層ニューラルネットワークは、予め蓄積されたエピソードデータを用いて訓練される。予め蓄積されたエピソードデータを用いてもステップ毎、エピソード毎に深層ニューラルネットワークを訓練できる。この訓練の際には強化学習で利用されるランダムな行動選択はしない。このエピソードデータを用いた深層ニューラルネットワークの訓練による学習をエピソードデータ学習と呼ぶ。

40

【0025】

エピソードデータ学習では学習パターンの設定に応じて、効率的に深層ニューラルネットワークを訓練でき、例えば、少ないエピソード数で学習が完了する、といった観点で効率化でき、学習結果の深層ニューラルネットワーク10によるシミュレーション結果は成功しやすくなる。

【0026】

行動制御計画装置は、個別のシミュレーション結果の強化学習の結果で得られたエピソードデータを利用し、複数のシミュレーション結果（例えば、拿捕、見失い、好位置など

50

)について予め学習パターンを設定して、学習パターンに一致するエピソードデータを予め蓄積されたエピソードデータから抽出してエピソードデータ学習を行う。

【0027】

なお、他移動体のシミュレーション結果は、上記の3種類に限定されるものではなく、予め設定された種類のシミュレーション結果を設定すればよい。シミュレーション結果は、自移動体と他移動体の位置関係を初期値から変化させて、所定の終了条件を満足した時点(エピソードの終了時点)の他移動体の位置から後述するように判定することができる。

【0028】

これにより、行動制御計画装置は、他移動体の複数の行動パターンに対して、状態量に対する行動(方策)の選択を深層ニューラルネットワークで学習する。すなわち、深層ニューラルネットワークは他移動体の行動パターンに応じた自移動体の行動制御計画を学習する。学習の完了は、複数のシミュレーション結果の全てに対して、エピソードの処理を試行し、各試行が成功することをもって判断する。

10

【0029】

以上により、シミュレーション結果毎の強化学習の結果であるエピソードデータを利用して、強化学習が成功するに至るエピソード毎の成功と失敗のステップ数を分析して学習パターンを設定する。

【0030】

そして、行動制御計画装置は、複数のシミュレーション結果についてエピソードデータ学習をすることにより、複数の行動パターンに対する自移動体の行動制御方法を計画できる。本実施例の行動制御計画装置は、従来の強化学習だけでは実現することが困難な、自移動体と他移動体の行動の複雑な関係において、深層ニューラルネットワークに効率よく行動制御を学習させることが可能となる。

20

【0031】

<処理の概要>

図1の行動制御計画装置の処理の一例を示す図により、本発明の一実施の形態に係わる、強化学習、エピソード分析、エピソードデータ学習による、行動制御計画の処理の概要を示す。まず、行動制御計画装置105で行われる処理の全体の流れを図2のフローチャートにより説明する。

【0032】

多様な環境や、他移動体の行動パターンに対して、自移動体の行動制御を計画するために、行動制御計画装置105は、まず、個別のシミュレーションの条件に対して自移動体の行動制御の強化学習を行う(201)。この強化学習においては、各移動体の状態に応じた深層ニューラルネットワーク10による行動選択を1ステップずつ進めて、自移動体の目的が達成されて成功、もしくは失敗するまでの1つのエピソードの処理において、状態量と方策、状態価値の関係を定める関数である深層ニューラルネットワーク10の訓練を行う。

30

【0033】

行動制御計画装置105は、エピソードの処理を繰り返して、必ず、もしくは十分に目的が達成される関数を得ることで学習を完了する。行動制御計画装置105は、学習が完了に至るまでのエピソードのステップ毎の状態量と、行動及び、報酬のデータを、強化学習過程の順序に対する全てのエピソードによる学習結果として、エピソードデータベース122に蓄積する(202)。

40

【0034】

行動制御計画装置105(図1参照)は様々な行動パターンに対する学習結果を蓄積する。また同一の学習パターンに対して、一つの学習結果だけでなく、多数の学習結果を蓄積する。これらの学習結果の一つのエピソードのデータをエピソードデータと呼ぶ。

【0035】

行動制御計画装置105は、エピソードデータベース122に予め蓄積されたエピソードデータに対して、個別の条件毎に強化学習を実施してエピソードデータを収集し、学習

50

過程の順にエピソードの種別（成功と失敗）、エピソードのステップ数を集計して、学習過程のエピソードを分析する（203）。

【0036】

これにより、行動制御計画装置105は、学習過程のエピソード順の、成功と失敗と、それらのステップ数を集計して学習パターンを設定する。

【0037】

そして、行動制御計画装置105は、複数の条件についての学習パターンの設定に従って、エピソードデータをエピソードデータベース122から収集する。そして、行動制御計画装置105は、学習過程の順にエピソードデータを用いて状態量と方策、状態値の関数（深層ニューラルネットワーク）の訓練を行う、エピソードデータ学習を行う（204）。学習過程の順のエピソード完了毎に、移動体の目的の達成（シミュレーション結果）を判定することで学習の完了を判定する。

10

【0038】

行動制御計画装置105は、エピソードデータ学習の学習結果を蓄積する（205）。エピソードデータ学習におけるエピソードデータは、既にエピソードデータベース122に蓄積されたものであり、エピソードデータそのものの蓄積は不要であり、学習過程で利用したエピソードデータのリストと、エピソードデータ学習における、エピソード毎の移動体の目的達成の判定結果（エピソードの種別）が蓄積されればよい。

【0039】

以上が、図2のフローチャートによる、強化学習、エピソード分析、エピソードデータ学習による、行動制御計画の処理の全体の流れの説明である。上記処理によって、強化学習で出力されたエピソードデータがエピソードデータベース122に蓄積され、エピソード分析で設定された学習パターンに応じたエピソードデータによって深層ニューラルネットワーク10が訓練され、訓練の結果はエピソードデータ学習で利用したエピソードデータのリストと、エピソード毎の移動体の目的達成の判定結果が蓄積されていく。

20

【0040】

図1の処理の一例を示す図は、本実施例の処理の概要において、必要となる処理機能と、処理の順序と、データの授受を矢印で示したものである。図中縦方向を強化学習部101、エピソード分析部102、エピソードデータ学習部103で分類し、図中横方向を処理の主体として、オペレータ104と、行動制御計画装置105で分類している。オペレータ104の処理とは利用者などのマニュアルによる処理を意味する。

30

【0041】

強化学習部101においては、自移動体が移動する空間の環境や、他移動体の行動パターン、また自移動体の移動方向、速度といった行動の設定などのシミュレーションを行う条件設定（シミュレーション条件設定部）111を処理する。強化学習部101は、オペレータ104から指定された条件（シミュレーション条件）で、シミュレータ112により他移動体と、自移動体の行動をシミュレーションする。

なお、シミュレータ112は、既に学習が完了した深層ニューラルネットワーク（機械学習モデル）を用いてシミュレーションを実施する。なお、シミュレータ112は、周知又は公知の技術を適用すればよく、移動体の種類に応じたシミュレータ112を使用することができる。また、移動体の種類は同一である必要はなく、例えば、自移動体が艦船、他移動体が潜水艦でもよいし、自移動体が航空機、他移動体が艦船であってもよい。

40

【0042】

シミュレータ112では、ステップ（所定の時間間隔）毎に状態更新113と行動選択114の処理がなされる。特に行動選択114では、強化学習のために、状態量と方策の関数から行動を選択するだけでなく、行動を乱択する処理も行う。

【0043】

自移動体の目標の成功、失敗までの一つのエピソードはステップ毎の処理を含み、ステップの状態量、行動、報酬のデータを用いて状態量と方策、行動値の関数である深層ニューラルネットワーク10の訓練を、学習器115で処理する。

50

【 0 0 4 4 】

シミュレーション結果のステップ毎の状態量、行動、報酬のデータは行動制御計画装置 1 0 5 のメモリ（記憶装置）に蓄積されており、訓練データ取得機能 1 1 6 により訓練用のデータがシミュレータ 1 1 2 から取得され、深層ニューラルネットワーク訓練 1 1 7 を処理する。自移動体の目的が達成されるように深層ニューラルネットワーク 1 0 の訓練が進み、強化学習が完了したら、強化学習部 1 0 1 における処理過程についてのエピソードデータが、エピソードデータベース（エピソードデータ蓄積部）1 2 2 に蓄積される。

【 0 0 4 5 】

エピソード分析部 1 0 2 においては、エピソードデータベース 1 2 2 から、エピソードの種別や他移動体のシミュレーション結果毎に、強化学習過程のエピソードデータを収集して、エピソード分析機能 1 2 1 により、学習パターンを生成する。

10

【 0 0 4 6 】

エピソードデータ学習部 1 0 3 では、対象とする複数の他移動体のシミュレーション結果について生成された学習パターンを用いて、エピソードデータ学習用に学習パターン設定 1 3 1 の処理を行う。

【 0 0 4 7 】

そして、学習パターン設定 1 3 1 で設定された学習パターンによりエピソード選択 1 3 2 の処理を行って、エピソードデータ学習 1 3 3 の処理を行う。エピソードデータ学習においても深層ニューラルネットワーク 1 0 の訓練のためには、学習器 1 1 5 を用いて処理を行う。

20

【 0 0 4 8 】

なお、条件設定 1 1 1 と学習パターン設定 1 3 1 はオペレータ 1 0 4 によるマニュアル設定を想定して図 1 に示したが、設定される条件を規則的に定めたり、学習パターンの設定方法が予め定められる場合には、行動制御計画装置 1 0 5 の処理としてもよい。

【 0 0 4 9 】

以上が図 1 の行動制御計画装置 1 0 5 で行われる処理の説明である。

【 0 0 5 0 】

< 行動制御計画 >

本発明の、移動体の行動制御計画の一例を示す。

【 0 0 5 1 】

図 3 のフローチャートにより、移動体の行動制御計画方法の一例を示す。フローチャートの処理は、強化学習部 1 0 1、エピソード分析部 1 0 2、エピソードデータ学習部 1 0 3 に分割して示している。各処理は行動制御計画装置 1 0 5 が実行するものとして記載する。

30

【 0 0 5 2 】

強化学習部 1 0 1 においては、まず、環境、他移動体の行動パターン、自移動体の行動等のシミュレーション条件を受け付けて、シミュレーション及び強化学習のためのパラメータについてのシミュレーション条件を取得する（3 0 1）。シミュレーション条件は、図 1 で示したようにオペレータ 1 0 4 が設定することができる。

【 0 0 5 3 】

状態量に対する方策と、行動価値の関数である深層ニューラルネットワーク 1 0 を訓練する処理では、強化学習のためのエピソードの繰り返しのステップ 3 0 2 からステップ 3 0 4 の範囲を繰り返して処理する。繰り返しの一回は一つのエピソードの処理のこととなる。

40

【 0 0 5 4 】

行動制御計画装置 1 0 5 は、繰り返しの 1 回について、1 エピソードの各ステップについて行動選択及び状態変化を深層ニューラルネットワーク 1 0 に学習させる（3 0 3）。1 つのエピソードの処理は、移動体の初期状態から、移動体の目的が達成されるか否かの成功、失敗までの、ステップ毎の処理である。

【 0 0 5 5 】

50

行動制御計画装置 105 は、ステップ毎、またエピソードの終了時点で、各ステップにおける状態量、行動、報酬のデータを用いて、深層ニューラルネットワーク 10 を訓練する。訓練時の行動選択では、行動は深層ニューラルネットワーク 10 による行動選択だけでなく、乱択（ランダムによる行動選択）も行う。

【0056】

また、行動制御計画装置 105 は、1つのエピソードの処理の終了時点で、訓練された深層ニューラルネットワーク 10 を用いた行動選択によるエピソードの処理をテスト（試行）し、所定の目的を達成したらエピソードの種別を成功として学習完了とする。この場合、例えば、エピソードが5回連続でテストに成功したら学習完了とする、といった学習完了の判定でもよい。

10

【0057】

行動制御計画装置 105 は、深層ニューラルネットワークの学習が完了したら、学習結果の深層ニューラルネットワークを用いてエピソードの試行（305）を行うことにより、学習が完了したこと、また移動体の行動内容を確認する。もしくは、学習が完了したことを確認できるようにログなどの記録を残すことができる。

【0058】

そして強化学習部 101 は、強化学習の結果（エピソード）のデータベース蓄積（306）で、強化学習での学習過程の全てのエピソードの、ステップ毎の状態量、行動、報酬のデータを、学習過程についてのエピソードデータとしてエピソードデータベース 122 に蓄積する。なお、強化学習部 101 は、エピソードデータを蓄積する際には、他移動体のシミュレーション結果とエピソードの種別を含めてエピソードデータベース 122 へ蓄積する。

20

【0059】

なお、このフローチャートにおいて、強化学習部 101 は、一通りのフローで示しているが、他移動体の行動パターンについての強化学習を、同一の行動パターンも含み、多数回実行する。

【0060】

エピソード分析部 102 においては、まず他移動体のシミュレーション結果毎に、エピソードデータベース 122 からエピソード収集（分析用）（307）を行う。この際には、強化学習における学習過程の順序についての情報も取得する。

30

【0061】

学習成功のエピソード失敗 / 成功過程収集（308）では、強化学習を実施した深層ニューラルネットワーク 10 がシミュレーションで成功した学習過程について、エピソード分析部 102 がエピソードの成功と失敗、そしてエピソードデータのステップ数を収集し、一覧化する。そして、学習過程におけるエピソードデータの順に、エピソード順の間隔毎に、失敗と成功の件数と各エピソードデータのステップ数を集計する。

【0062】

行動制御計画装置 105 は、エピソードデータの集計結果より、エピソード順の成功と失敗のサンプル（エピソードデータ）の件数と、ステップ数の範囲を設定して、学習パターンを構成する（309）。

40

【0063】

なお、このフローチャートにおいてエピソード分析部 102 は、一通りのフローで示しているが、他移動体の異なる行動パターンについても同様にエピソード分析を行う。

【0064】

エピソードデータ学習部 103 においては、まずステップ 301 と同様の条件と併せて、学習パターンを含めて、条件（エピソード条件）を取得（学習パターンを含む）する（310）。なお、この条件（エピソード条件）には、複数の他移動体のシミュレーション結果の指定も含むことができる。

【0065】

そしてエピソード学習のために、学習パターンで指定されている学習過程の順序におけ

50

るエピソードの種別（成功と失敗）、及びステップ数に該当するエピソードデータをエピソードデータベース 1 2 2 から検索して取得する、エピソード収集（学習用）を行う（3 1 1）。

【0 0 6 6】

エピソードデータ学習部 1 0 3 において、状態量に対する方策と、行動価値の関数である深層ニューラルネットワークを訓練する処理は、収集したエピソードについて、ステップ 3 1 2 からステップ 3 1 4 の範囲でエピソード毎の繰り返しで行う。エピソード毎の繰り返しは、予め設定された終了条件（例えば、処理済みのエピソード数等）が成立するまで行うことができる。

【0 0 6 7】

行動制御計画装置 1 0 5 は、繰り返しの 1 回について、1 エピソードの行動選択及び状態データを学習する（3 1 3）。ここでの学習は、行動を乱択することもある、いわゆる強化学習ではなく、エピソードデータを直接に用いた学習である。1 つのエピソードデータは、移動体の初期状態から、エピソードが成功又は失敗するまでの、ステップ毎の処理である。

【0 0 6 8】

エピソードデータ学習部 1 0 3 は、エピソードデータのステップ毎に処理を進めるが、シミュレータ 1 1 2 を利用して移動体の行動を選択したり、状態を更新して、状態量、報酬を計算することはない。

【0 0 6 9】

エピソードデータ学習部 1 0 3 は、ステップ毎にエピソードデータの状態量、行動、報酬を各移動体に設定していく。このデータをメモリに格納し、ステップ毎、またエピソードの終了時点で、状態量、行動、報酬のデータを用いて、深層ニューラルネットワークを訓練する。

【0 0 7 0】

エピソードデータ学習部 1 0 3 は 1 つのエピソードの処理の終了時点で、訓練された深層ニューラルネットワークを用いた行動選択によるエピソードの処理をテストする。この際に、全ての他移動体のシミュレーション結果に対するエピソードの処理をテストし、全て成功したら学習完了とする。また各シミュレーション結果についてエピソードが 5 回連続で成功したら学習完了とする、といった終了条件による判定をしてもよい。

【0 0 7 1】

学習が完了したら、エピソードデータ学習部 1 0 3 は、学習結果でのエピソード試行（3 1 5）を行うことにより、学習が完了したこと、また各移動体の行動内容を確認する。もしくは確認できるようにログなどの記録を出力する。エピソードデータ学習部 1 0 3 は、学習で指定された、他移動体のシミュレーション結果の全てに対して試行する。

【0 0 7 2】

エピソードデータ学習部 1 0 3 は、エピソードデータ学習結果のデータベース蓄積（3 1 6）を行う。なお、学習に利用したエピソードデータの登録はしない。エピソードデータ学習の学習過程、例えば学習完了までにかかったエピソード数や利用したエピソードデータなどをエピソードデータベース 1 2 2 に蓄積する。

【0 0 7 3】

以上が、図 3 のフローチャートによる、移動体の行動制御計画方法の一例の説明である。上記処理によって、深層ニューラルネットワーク 1 0 が出力したエピソードデータが強化学習の結果としてエピソードデータベース 1 2 2 に格納され、エピソード分析部 1 0 2 は、エピソードの成功と失敗及びエピソードデータのステップ数を収集し、エピソードの成功と失敗の件数と、ステップ数の範囲を設定して学習パターンを後述するように生成する。

【0 0 7 4】

エピソードデータ学習部 1 0 3 は、学習パターンに対応するエピソードデータをエピソードデータベース 1 2 3 から選択して学習器 1 1 5 を用いてエピソードデータ学習 1 3 3

10

20

30

40

50

を実施して深層ニューラルネットワーク 10 訓練を実施する。そして、エピソードデータ学習 133 は、訓練された深層ニューラルネットワーク 10 について試行を実施してからエピソードデータベース 122 にエピソード数等の学習過程のデータを蓄積する。

【0075】

< 行動制御計画装置の構成 >

次に、図 4 の行動制御計画装置 105 の機能の一例を説明する。図 4 は、図 3 のフローチャートで説明した内容を実現する行動制御計画装置 105 の機能である。

【0076】

行動制御計画装置 105 は、プロセッサ 11 と、メモリ 12 と、ストレージ装置 13 と、入出力装置 14 を含む計算機で構成される。なお、入出力装置 14 は、キーボードやマウス又はタッチパネル等の入力装置と、ディスプレイなどの出力装置を含む。

10

【0077】

行動制御計画装置 105 のメモリ 12 には、シミュレータ 400 と、強化学習 / エピソードデータ学習部 410 と、エピソード分析部 430 の機能部がプログラムとしてロードされて、プロセッサ 11 によって実行される。なお、各プログラムが利用するデータ（行動環境データベース 418、エピソードデータベース 419 等）はストレージ装置 13 に格納してもよい。

【0078】

強化学習 / エピソードデータ学習部 410 は、深層ニューラルネットワーク（DNN：Deep Neural Network）10 の訓練に関する処理が、強化学習部（図 1 の 101 に相当）とエピソードデータ学習部（図 1 の 103 に相当）と共通であるため、強化学習部とエピソードデータ学習部を合わせて一つの機能ブロック（410）とした。しかし、エピソードデータ学習部は強化学習とは異なる学習方法であるため、機能ブロック（410）の内側に、エピソードデータ学習部 420 の機能ブロックを配置した。

20

【0079】

シミュレータ 400（図 1 の 112 に相当）の機能は強化学習 / エピソードデータ学習部 410 で利用される。シミュレータ 400 は、図 3 における 1 エピソードの行動選択・状態変化の学習（303）の処理において、ステップ毎の処理を進めるための機能である。

【0080】

シミュレータ 400 は行動環境設定機能 401、移動体行動処理機能 402、移動体状態量算出機能 403、移動体状態判定機能 404 を有する。

30

【0081】

行動環境設定機能 401 は、自移動体、他移動体の初期位置や障害物などの環境を設定する機能であり、本機能により初期状態を決める。

【0082】

移動体行動処理機能 402 により、ステップ毎で自移動体の行動を選択し、移動処理を行う。移動体行動処理機能 402 他移動体の移動も行い、これによりステップ毎の状態が更新されることとなる。移動体行動処理機能 402 は、図 1 の状態更新 113 と、行動選択 114 に対応する。

【0083】

移動体状態量算出機能 403 により、ステップ毎で状態が更新された際に、状態量を計算する。状態量とは、自移動体が予め設定された目的を達成するために必要となる情報であって、例えば、他移動体との距離や、障害物との距離、他移動体の方位、他移動体の速度、などである。

40

【0084】

移動体状態判定機能 404 により、計算された状態量から自移動体のそのステップにおける状態を判定する。状態は、例えば自移動体と他移動体の距離によれば、他移動体を検知できている良好な位置にある状態（好位置）、他移動体を検知できない距離となって見失った状態、検知できているが近くなりすぎて他移動体に見つかり拿捕される状態、といった記号的な状態名称の判定が行われる。なお、本実施例の自移動体は、艦船の例を示す。

50

【 0 0 8 5 】

移動体状態量算出機能 4 0 3 は、この判定結果を利用して、他移動体の見失い、自移動体の拿捕といった状態ならばエピソードは失敗、また所定ステップ数を満たして良好な位置（好位置）ならば成功と判定しエピソードを終了する。また、報酬の値は、この状態に基づいて設定される。

【 0 0 8 6 】

以上の機能により、図 3 のステップ 3 0 3 の処理（行動選択及び状態変化の学習）が実現される。なおシミュレーションを実行することは強化学習では必要だが、エピソードデータ学習そのものには不要な機能である。しかし、一般にシミュレータ 4 0 0 は自移動体、他移動体、環境に関するデータを保持し、エピソードデータにアクセスする機能をプログラムとして有しており、当該プログラムの機能はエピソードデータ学習でも利用され得る。

10

【 0 0 8 7 】

強化学習 / エピソードデータ学習部 4 1 0 の枠内は学習を実現するための各機能を示している。このうち、特に強化学習に係わる機能について示す。強化学習 / エピソードデータ学習部 4 1 0 は、方策・状態価値関数構成機能 4 1 1、エピソード処理機能 4 1 2、報酬設定機能 4 1 3、DNN 訓練機能 4 1 4、DNN 訓練用データサンプリング機能 4 1 5、方策・状態価値算出機能 4 1 6、エピソードデータ蓄積機能 4 1 7、行動環境データベース 4 1 8、エピソードデータベース 4 1 9 の構成要素を有する。

【 0 0 8 8 】

方策・状態価値関数構成機能 4 1 1 により、強化学習の開始前に、状態量を入力し、選択される行動、すなわち方策と、状態価値を出力とする深層ニューラルネットワーク 1 0 を構成する。

20

【 0 0 8 9 】

エピソード処理機能 4 1 2 により、1 つのエピソードに対して、初期状態からステップ毎に処理を進めて、移動体が目的を達成するか否かによるエピソードの終了までの処理を進める。エピソード処理機能 4 1 2 は、図 3 のステップ 3 0 3 の処理を行う機能であり、シミュレータ 4 0 0 の各機能、及び報酬設定機能 4 1 3、DNN 訓練機能 4 1 4、DNN 訓練用データサンプリング機能 4 1 5、方策・状態価値算出機能 4 1 6 の呼び出しなど、各種制御を行う機能である。

30

【 0 0 9 0 】

報酬設定機能 4 1 3 は、移動体状態判定機能 4 0 4 の結果で得られた状態に応じて報酬を決める機能である。

【 0 0 9 1 】

DNN 訓練機能 4 1 4 により、深層ニューラルネットワーク 1 0 のパラメータを、ステップ毎の状態量、行動、報酬のデータを用いて最適化することで、深層ニューラルネットワーク 1 0 を訓練する。

【 0 0 9 2 】

DNN 訓練用データサンプリング機能 4 1 5 により、深層ニューラルネットワーク 1 0 を訓練するためのデータを、エピソードにおけるステップ毎のデータから特に最適化を効率的に進めるために、サンプリングする。エピソードについては一つのエピソードだけでなく、複数のエピソードに亘るサンプリングとしてもよい。

40

【 0 0 9 3 】

方策・状態価値算出機能 4 1 6 は、深層ニューラルネットワーク 1 0 の入力となる状態量を与えることで、深層ニューラルネットワーク 1 0 の出力である方策と状態価値を取得する機能である。方策を取得するとは、行動を選択するということであり、移動体行動処理機能 4 0 2 でこの方策・状態価値算出機能 4 1 6 を呼び出すことで行動を決定する。

【 0 0 9 4 】

エピソードデータ蓄積機能 4 1 7 により、図 3 のステップ 3 0 6 におけるエピソードデータをエピソードデータベース 4 1 9 に蓄積する。

50

【0095】

行動環境データベース418は、自移動体の行動環境のデータが格納されるデータベースである。

【0096】

エピソードデータベース419は、エピソードデータが格納されるデータベースで、図1のエピソードデータベース122に相当する。

【0097】

エピソードデータ学習部420の各機能について説明する前に、図3のフローチャートの順に合わせ、エピソード分析部430の各機能について説明する。

【0098】

エピソード分析部430は、エピソードデータ収集機能431と、エピソードデータ描画機能432、ステップ毎エピソード集計機能433、学習パターン構成機能434で構成される。なお、エピソード分析部430は、図1のエピソード分析部102に相当する。

【0099】

エピソードデータ収集機能431により、他移動体のシミュレーション結果毎に、強化学習の学習過程における、エピソードデータの収集を行う。

【0100】

エピソードデータ描画機能432により、エピソードにおける環境、自移動体、他移動体のステップ毎の、例えば配置もしくはマップ(地図)を描画する。またステップ毎の状態量の変化などの描画も含む。マニュアルでのオペレータ104によるエピソードの内容確認のための機能である。なお、本実施例では、画像データの生成を描画とする。

【0101】

ステップ毎エピソード集計機能433により、学習過程における、エピソードの成功と失敗、ステップ数を一覧化し、エピソードの順に成功と失敗の件数とエピソードデータのステップ数を集計する。

【0102】

学習パターン構成機能434により、エピソード順の間隔毎に、エピソードの種別(成功と失敗)の件数と、エピソードデータのステップ数の範囲を集計することで、学習パターンを構成する。

【0103】

強化学習/エピソードデータ学習部410でエピソードデータ学習部420を構成する各機能について説明する。エピソードデータ学習部420は、学習パターン取得機能421、エピソードデータ収集・一覧化機能422、エピソードデータ学習機能423で構成される。

【0104】

学習パターン取得機能421により、他移動体のシミュレーション結果についての学習パターンを取得する。図3のステップ310で利用される機能であり、この際、行動環境設定機能401により、各種の条件も取得する。なお、学習パターン取得機能421が取得する学習パターンは、エピソード分析部430の学習パターン構成機能434が出力した学習パターンである。

【0105】

エピソードデータ収集・一覧化機能422により、学習パターンで指定されたエピソードの条件(エピソード条件)に該当するエピソードデータをエピソードデータベース419から収集する。エピソードデータ収集・一覧化機能422は収集した結果をエピソードデータ学習で利用されるデータとして一覧化する。

【0106】

エピソードデータ学習機能423により、エピソードデータ学習の処理を行う。エピソードデータ学習機能423は、強化学習におけるエピソード処理機能412に相当し、各機能を制御する機能を有する。エピソードデータ学習の開始前には、方策・状態価値関数構成機能411により深層ニューラルネットワーク10を構成する。

10

20

30

40

50

【0107】

強化学習/エピソードデータ学習部410は、エピソードデータのステップ毎のデータをメモリ12に格納し、DNN訓練機能414、DNN訓練用データサンプリング機能415を利用して、深層ニューラルネットワーク10の訓練を行う。またエピソードデータ学習が完了したら、エピソードデータ蓄積機能417により、学習結果をエピソードデータベース419に格納する。

【0108】

上記各機能はプログラムとしてメモリ12にロードされてプロセッサ11によって実行される。プロセッサ11は、各機能部のプログラムに従って処理を実行することによって、所定の機能を提供する機能部として稼働する。例えば、プロセッサ11は、行動環境設定プログラムに従って処理を実行することで行動環境設定機能401を提供する。他のプログラムについても同様である。さらに、プロセッサ11は、各プログラムが実行する複数の処理のそれぞれの機能を提供する機能部としても稼働する。計算機及び計算機システムは、これらの機能部を含む装置及びシステムである。

10

【0109】

以上が、図4の構成による、移動体の行動制御計画装置105が有する機能の一例の説明である。

【0110】

<移動体の行動>

本発明における移動体の行動の目的と環境の例を示す。

20

【0111】

自移動体の行動の目的は、他移動体の追跡(追尾)、又は設定された領域への他移動体の侵入を検知する哨戒とする。

【0112】

図5に、追跡における自移動体と、他移動体の行動環境の例を示す。

【0113】

枠501において、2次元空間である平面の座標をとるためのX座標軸502とY座標軸503が設定される。なお行動領域は平面を示したが、深さ方向を含む3次元空間としてもよい。

【0114】

自移動体504は、直進、旋回、また変速して空間を自由に移動できる。初期位置と初期速度は予めシミュレーション条件で設定されるが、方向、速度はステップ毎の行動選択により変化する。

30

【0115】

他移動体505は、初期位置から、一定の速度で、一方向506に移動する。初期位置、速度、方向は予め設定され、この組み合わせで一つの他移動体505の行動パターンとなる。

【0116】

自移動体504には検知距離507が設定され、検知距離507を半径とする円508の内側の他移動体505を検知する。追跡開始の状態、すなわち初期位置では他移動体505が自移動体504に検知されるように設定される。

40

【0117】

円508の外側は自移動体504が他移動体505を検知できないので見失う範囲509となる。自移動体504が他移動体505を見失う範囲509に入った場合、自移動体504は見失いというシミュレーション結果となり、他移動体505の追跡という目的は達成されず、エピソードの種別は失敗となる。

【0118】

また自移動体504は、他移動体505に接近しすぎた場合には、他移動体505によって検知され、捕まるとする。この状態を拿捕として、追跡の継続という目的は達成されず、エピソードの種別は失敗となる。このために、自移動体504には拿捕範囲510が

50

設定され、円形の拿捕範囲 5 1 0 の内側に他移動体 5 0 5 が入ることで、拿捕を判定する。

【 0 1 1 9 】

シミュレーションの開始から所定の、もしくは設定されたステップ数を経過するまで、見失い、及び拿捕の状態になることがなかった場合、自移動体 5 0 4 の追跡の目的は達成されたと判定され、エピソードの種別は成功となる。

【 0 1 2 0 】

以上が図 5 に示した、追跡における行動環境の例の説明である。

【 0 1 2 1 】

図 6 に、哨戒における自移動体、他移動体の行動環境の例を示す。

【 0 1 2 2 】

哨戒領域 6 0 1 は予め設定される。自移動体 6 0 2 は哨戒領域 6 0 1 内を初期状態（位置、速度）から自由な方向 6 0 5 に移動できる。自移動体 6 0 2 の移動方向、速度はステップ毎の行動選択により変化する。

【 0 1 2 3 】

自移動体 6 0 2 は他移動体を検知する能力を有し、検知距離 6 0 3 を半径とする円 6 0 4 の範囲に他移動体があれば、他移動体を検知したと判定する。

【 0 1 2 4 】

他移動体は、図 6 の例では、垂直方向で哨戒領域 6 0 1 に侵入する他移動体 6 0 6 と、水平方向で侵入する他移動体 6 0 7 といったように行動にバリエーションがある。侵入位置は様々、もしくはランダムとする。侵入してきた他移動体 6 0 6、6 0 7 に対して、自移動体 6 0 2 が検知距離 6 0 3 に入れば検知に成功したと判定する。

【 0 1 2 5 】

他移動体が検知されることなく哨戒領域 6 0 1 を通過したら検知に失敗したとする。自移動体 6 0 2 が他移動体を検知することが目的なので、検知成功ならプラスの報酬が付与され、検知失敗ならマイナスの報酬が与えられる。

【 0 1 2 6 】

エピソードとしては、一つの他移動体の行動に対して検知の成功、もしくは失敗をもって、エピソードの種別を成功又は失敗に分類することができる。またエピソードのステップ数を予め定めておいて、また他移動体の数は多数として、強化学習してもよい。この場合、多くの報酬が得られるような自移動体 6 0 2 の行動が得られることとなる。

【 0 1 2 7 】

他移動体の行動は、発生位置をランダムとすれば自移動体 6 0 2 の普遍的な哨戒行動が得られることとなるが、予め発生位置が指定されれば他移動体の特定の侵入の傾向に対する自移動体 6 0 2 の哨戒行動が得られることとなる。

【 0 1 2 8 】

哨戒行動は、哨戒領域 6 0 1 の大きさに依存し、また他移動体の哨戒領域 6 0 1 への侵入行動パターンにも依存し、環境についての条件となる。

【 0 1 2 9 】

図 7 A に哨戒領域と座標系の例を示す。図 7 B に哨戒領域の構成と他移動体の行動パターンの例を示す。図 7 A では、哨戒領域 7 0 1 を四角領域とし、中心（図心、重心）を原点として、X 軸 7 0 2、Y 軸 7 0 3 をとる。X 方向の範囲 L_x 7 0 4、Y 方向の範囲 L_y 7 0 5 により領域の大きさを設定する。範囲 L_x 7 0 4、 L_y 7 0 5 の半分を $h L_x$ 7 0 6、 $h L_y$ 7 0 7 とすれば、座標値としては $\pm h L_x$ 、 $\pm h L_y$ となる。

【 0 1 3 0 】

図 7 B に他移動体の行動パターンの例を示す。哨戒領域 7 0 1 に対して外側にオフセットを付加して広げた領域を他移動体発生領域 7 1 0 とする。他移動体は哨戒領域 7 0 1 の外部、かつ他移動体発生領域 7 1 0 の内部で発生する。

【 0 1 3 1 】

哨戒領域 7 0 1 は四角形であり、他移動体は領域の 4 辺の上 7 1 1、下 7 1 2、右 7 1 3、左 7 1 4 のそれぞれから移動方向 7 1 5、7 1 6、7 1 7、7 1 8 で侵入する。

10

20

30

40

50

【0132】

他移動体の哨戒領域701への侵入件数としては、予め設定されたステップ毎に1つずつ侵入させる、哨戒領域701内の他移動体の数を一定とする、などと決めることができる。

【0133】

以上が、図6、図7A、図7Bに示した、哨戒における自移動体、他移動体の行動環境の説明である。

【0134】

以上が、本発明における移動体の行動の目的と環境についての説明である。

【0135】

<強化学習>

本発明における強化学習の内容について示す。強化学習は、図1、図3における強化学習部101及び図4の強化学習/エピソードデータ学習部410の処理に相当する。

【0136】

強化学習は、一般的にはマルコフ決定過程MDP(Markov Decision Process)のモデルに基づいている。マルコフ決定過程は、ある状態において行動を選択して状態を更新し、この際の行動もしくは状態には報酬として表現される価値(又は状態価値)が伴う、というモデルである。

【0137】

強化学習の方式には、行動選択を行動価値のテーブル値により行うQ学習、また行動価値関数を深層ニューラルネットワークとして行動価値から行動を選択するDQN(Deep Q-learning Network)、方策(行動選択)関数と状態価値関数を深層ニューラルネットワークとしたA3C(Asynchronous Advantage Actor-Critic)といった方法がある。

【0138】

深層ニューラルネットワークを利用したアルゴリズムは他にも存在するが、いずれも強化学習として、状態、行動、報酬によってモデルが訓練される。本発明ではA3Cを対象に説明する。

【0139】

強化学習は、エピソードの繰り返しにおいて、エージェントが状態に対して最適な行動を選択するように方策関数、状態価値関数を最適化する方法である。エージェントは学習の主体であって、例えば、本発明では自移動体のことである。エージェントは一つのエピソードにおいてステップ毎に、状態に対して行動を選択して、次の状態を得る。

【0140】

強化学習では、行動選択の結果や状態に対して報酬を与える。エピソードの終了条件を満たすならエピソードを終了する。エピソードを繰り返して、状態に対して良好な状態が得られるように正しく行動選択できるようになれば強化学習は終了となる。

【0141】

A3CはAsynchronous Advantage Actor-Criticの略であり、非同期処理、アドバンテージによる評価、アクター-クリティックとしての方策、状態価値関数モデルが特徴である。非同期処理は個別に用意された環境のそれぞれでエージェントのエピソード処理を行う。

【0142】

この場合、それぞれのエージェントが単一方策関数と、状態価値関数を訓練することで、学習を効率化する工夫である。アドバンテージとは行動価値と状態価値の差であって、行動選択の良さを表す量である。アクター-クリティックとは、行動選択はアクターとして、状態価値はクリティックとして、別々のモデルで計算する方式であることを意味する。

【0143】

本発明での方策関数、状態価値関数の深層ニューラルネットワークによる構成の例を図

10

20

30

40

50

8に示す。深層ニューラルネットワークは状態量801を入力とし、行動802と、状態価値803を出力とする。

【0144】

状態量801から行動802までのネットワークが方策関数、状態量801から状態価値841までのネットワークが状態価値関数である。状態量801は、自移動体から他移動体までの距離821と、自移動体から見た他移動体の方位822と、自移動体の速度823である。

【0145】

行動802は自移動体の旋回行動の右831、左832、前833、変速行動の減速834、加速835である。状態価値803はそれ自体が一つの数値として状態価値841のみが項目となる。

10

【0146】

深層ニューラルネットワークは入力層811と、中間層812と、出力層813の多層構造として構成され、特に中間層812は複数の層となっている。なお図8ではノード(節点)を白丸で描いているが、全てのノードと層は描かずに点線で結合関係を表し、点線丸囲いでノードの集まりを表すこととして描いた。

【0147】

なお各ノードは関数であって、ノード同士を結ぶ線が入出力の関係を表す。状態量801は直接に入力層811に対応する。出力層813の各ノードは行動802、状態価値803の各項目に対応する。

20

【0148】

行動802の各項目に対応する出力層813のノードは総和が1となるように、確率の意味で事象分けされるようにSoftmax関数が活性化関数として定義される。状態価値841に対応する出力関数にはLinear関数で定義する。なお中間層812の各ノードにはLer関数、Sigmoid関数などが定義される。

【0149】

エピソードにおいてステップ毎のデータはメモリ12に格納され、例えば10ステップに1回といった、ステップ毎の所定のタイミングや、エピソード終了時に深層ニューラルネットワークが訓練される。各ステップのデータは状態量、行動、報酬である。

【0150】

図9A、図9Bにエピソードにおける状態量、行動、報酬の一連のステップデータ(エピソードデータ)の構成を示す。図9Aは失敗した場合のステップデータ900Aの例を示し、図9Bは成功した場合のステップデータ900Bの例である。

30

【0151】

図9A、図9Bの各図はそれぞれが1つのエピソードを構成するエピソードデータである。各図の項目は大きく分けるとステップ901と、状態量902と、行動903、報酬904、状態905である。状態905は深層ニューラルネットワークの訓練に利用しないが、参照、またエピソード分析で必要となる。

【0152】

状態量902は距離911、方位912、速度913であり、その数値は0以上1以下の範囲に正規化される。行動903は前914、右915、左916、減速917、加速918であり、一つのステップで選択される行動は一つであるため、一つの項目を1、他を0とする。なお、距離911は、自移動体から他移動体までの距離を示し、方位912は、自移動体から見た他移動体の方位を示し、速度913は、自移動体の速度を示す。

40

【0153】

報酬904はエピソードが失敗したら-1、成功したら1、それ以外のステップでは0とする。状態905は好位置、見失い、拿捕の3種類である。

【0154】

方策関数と状態価値関数の訓練の処理内容について説明する。深層ニューラルネットワークは状態量 s を入力として、行動の出力である方策を (s) 、状態価値の出力を $V($

50

s)とするネットワークである。

【0155】

ネットワーク中のノードに対して定義される活性化関数のパラメータをとする。方策は確率方策とも呼ばれ、状態sにおいて行動aをとる確率(a|s)とも記載する。方策関数、状態価値関数の訓練とは状態s、行動a、次状態s'、報酬rのデータから、それらの関係を予測する関係が得られるように、深層ニューラルネットワークのパラメータを最適化することである。

【0156】

状態価値関数V(s)は、方策π(s)の下で以下の式(1)で表現される。

【0157】

【数1】

$$V(s) = E_{\pi(s)}[r + \gamma V(s')] \quad (1)$$

【0158】

ここでEは、添え字である方策π(s)における期待値を意味する。割引率γは、次の状態価値(将来の価値)を現在の値に補正するための係数である。割引率γの値は、一例としては0.99であるが、強化学習のための設定として調整できる。方策の価値は状態sの分布π(s)における期待値として以下の式(2)で表される。

【0159】

【数2】

$$J(\pi) = E_{\rho^s}[V(s)] \quad (2)$$

【0160】

式(2)を割引済み報酬関数と呼ぶ。割引済み報酬関数の深層ニューラルネットワークのパラメータに関する変化には、式(3)で表現される方策勾配定理がある。

【0161】

【数3】

$$\nabla_{\theta} J(\pi) = E_{s \sim \rho^{\pi}, a \sim \pi(s)} [A(s, a) \nabla_{\theta} \log \pi(s, a)] \quad (3)$$

【0162】

ここで、∇は勾配(基底についての1階の偏微分)を意味する。期待値は状態sが方策π(s)の分布π(s)の巨る範囲、行動aは状態sに対する方策π(s)に巨る範囲による。行動価値関数Q(s, a)、アドバンテージ関数A(s, a)は以下であり、データで計算可能である。

【0163】

【数4】

$$Q(s, a) = r + \gamma V(s') \quad (4)$$

$$A(s, a) = Q(s, a) - V(s) \quad (5)$$

【0164】

報酬関数を最大とするようにネットワークパラメータθを最適化すれば、状態に対して良い報酬が得られるような方策と状態価値の関係が得られる。報酬関数の負(マイナス)を方策損失とすれば、方策損失を最小化する。

10

20

30

40

50

【 0 1 6 5 】

方策損失の他にも、行動価値と状態価値は一致していることが望ましいのでアドバンテージの絶対値の大きさを意味する価値損失も存在する。また状態に対して方策は一つに決まる方が望ましく、方策は確率的に決まるのでエントロピーでモデル化した正則化項も最適化に利用する。そこで損失関数を方策損失 L 、価値損失 L_v 、正則化項 L_{reg} により定義し、損失関数を最小化する。

【 0 1 6 6 】

【数 5】

$$L = L_{\pi} + c_v L_v + c_{reg} L_{reg} \quad (6)$$

10

【 0 1 6 7 】

ここで c_v 、 c_{reg} は係数である。

【 0 1 6 8 】

方策損失は、割引済み報酬関数の定義から以下となる。ここで、訓練に使う n はステップのデータ数である。

【 0 1 6 9 】

【数 6】

$$L_{\pi} = -J(\pi) = -\frac{1}{n} \sum_{i=1}^n A(s_i, a_i) \log \pi(a_i | s_i) \quad (7)$$

20

【 0 1 7 0 】

価値損失は、アドバンテージ関数の二乗とする。

【 0 1 7 1 】

【数 7】

$$L_v = \frac{1}{n} \sum_{i=1}^n A(s_0, a_0)_i^2 \quad (8)$$

30

【 0 1 7 2 】

正則化項はエントロピー $H(\pi(s))$ を計算して得る。

【 0 1 7 3 】

【数 8】

$$H(\pi(s)) = -\sum_{k=1}^{n_{action}} \pi(s)_k \log \pi(s)_k \quad (9)$$

40

$$L_{reg} = \sum_{i=1}^n H(\pi(s_i)) \quad (10)$$

【 0 1 7 4 】

ここで n_{action} は行動の数である。

【 0 1 7 5 】

最適化計算によるニューラルネットワークの訓練には勾配報 (gradient) を利用する。この訓練自体は深層学習における深層ニューラルネットワークの訓練と同じであ

50

る。

【 0 1 7 6 】

以上が方策関数と状態価値関数の訓練の処理内容についての説明である。

【 0 1 7 7 】

深層ニューラルネットワークの訓練を効率化するために、ステップ毎のデータのサンプリング方法について説明する。

【 0 1 7 8 】

一般的にデータを用いた最適化計算においては、サンプルによる目的関数値の値が、最適化の目的から外れた値である方が、モデルのパラメータの調整幅が大きくなり、最適化の効率がよい。

【 0 1 7 9 】

例えば、目的関数を最小化する場合、目的関数値が大きくなるデータを用いて最適化計算をすることで、目的関数の値を大幅に小さくするようにパラメータの値を調整できる。目的関数値の変化を大きくできることで、パラメータの調整幅を大きくできるのである。

【 0 1 8 0 】

目的関数は式(6)で表現される損失関数であり、この最小化が最適化の目的である。その項となっている方策関数、価値関数、正則化項は式(7)、(8)、(9)に示されるように、状態量 s と行動 a のデータを利用して計算されることを示す。これらの計算には方策関数 (s) もしくは $(a | s)$ と、アドバンテージ関数 $A(s, a)$ が利用される。

【 0 1 8 1 】

方策関数は深層ニューラルネットワークで状態量から行動を選択する関数であり、状態量 s から行動別の値を取得できる。アドバンテージ関数は式(5)で定義され、また式(4)の行動価値関数 $Q(s, a)$ を利用する。

【 0 1 8 2 】

行動価値関数 $Q(s, a)$ は式(4)から状態価値関数 $V(s')$ と報酬 r で計算される。状態価値関数は深層ニューラルネットワークで状態量から状態価値を算出する関数であり、状態量 s で値を取得できる。

【 0 1 8 3 】

行動価値関数においては、次状態 s' はエピソードデータでのステップで、次の行のステップの状態量 s を用いて計算する。最終行のステップについては状態価値関数 $V(s')$ をゼロとすればよい。以上から損失関数値は、ステップ別に状態量 s 、行動 a 、報酬 r から計算できる。

【 0 1 8 4 】

個々のステップに対して損失関数値が得られるので、損失関数値が大きくなる順にステップデータを並べて、大きな値となったステップの状態量、行動、報酬データをサンプリングすることで、訓練を効率化できる訓練用のデータが得られる。

【 0 1 8 5 】

多くのステップ毎のデータから訓練に必要な数のデータを損失関数値の大きさの順にサンプリングしても、また限られたステップデータから重複も許して確率的にサンプリングしてもよい。また損失関数値が小さいことも、それまでに訓練されてきた深層ニューラルネットワークのパラメータの特徴を反映しているので損失関数値が小さくなるステップデータも混ぜてサンプリングしてもよい。

【 0 1 8 6 】

以上が深層ニューラルネットワークの訓練を効率化するための、ステップデータのサンプリング方法についての説明である。処理の主体は行動制御計画装置 105 である。

【 0 1 8 7 】

一連の強化学習の処理手順の一例を、図10のフローチャートにより説明する。

【 0 1 8 8 】

ステップ 1001 からステップ 1018 はエピソードの繰り返しである。

10

20

30

40

50

【0189】

エピソードの処理においては、まず環境、自移動体、他移動体の状態をリセットする1002。

【0190】

ステップ1003からステップ1014は、1つのエピソードにおけるステップの繰り返しである。

【0191】

強化学習/エピソードデータ学習部410は、自移動体の行動を選択する(1004)。この処理は、方策関数に現在の状態量を入力することで行われる。

【0192】

そして強化学習/エピソードデータ学習部410は選択された行動を実行し、状態量を更新する(1005)。強化学習/エピソードデータ学習部410は、シミュレータ400の機能を利用して状態量を算出する。

【0193】

強化学習/エピソードデータ学習部410は、自移動体の状態量に基づき、エピソードの失敗に関するNG判定を行う(1006)。追跡の例ならば、見失い、拿捕の判定が失敗に相当する。

【0194】

強化学習/エピソードデータ学習部410は、NG判定1006がyes(NGである)ならば、エピソード失敗のフラグを設定する(1007)。

【0195】

次に、強化学習/エピソードデータ学習部410は、エピソードで実施したステップ数の判定を行う(1008)。この判定は、実施したステップ数が所定数を超えていれば、エピソードの終了と判定してステップ1009へ進み、そうでない場合にはステップ1011に進む。

なお、上記ステップ1007で、既に失敗フラグが設定されている場合は、ステップ1008の判定自体が不要であり、判定結果としてはno(終了ステップ数ではない)となる。またその場合にはエピソード終了フラグも設定する。

【0196】

ステップ数の判定(1008)でyes(処理したステップ数が所定の終了ステップ数である)の場合、強化学習/エピソードデータ学習部410は、成功フラグを設定する(1009)。これは、失敗ではなく、指定されたステップ数の間、処理を継続できたのでエピソードの種別が成功ということの意味する。処理が追跡の例ならば、好位置の状態を継続してきたことが成功に相当する。

【0197】

なお、哨戒の例については、NG判定(1006)でOK(no)の場合には、他移動体の検知として、検知数をカウントする。NG(yes)の場合には検知失敗の見失いと判定して見失い数をカウントする。いずれの場合も他移動体を環境から除外する処理を行う。ステップ数判定(1008)は、終了ステップ数だけの判定になり、特に成功フラグ設定(1009)は不要である。

【0198】

ステップ数判定(1008)でyesの場合には、強化学習/エピソードデータ学習部410がエピソード終了フラグを設定する(1010)。

【0199】

強化学習/エピソードデータ学習部410は、深層ニューラルネットワークの訓練のため、訓練データを設定する(1011)。このステップ1011は、訓練のためにメモリ12から一纏まりのデータを準備することであって、訓練を効率化するためのサンプリングといった処理も含まれる。

【0200】

そして、強化学習/エピソードデータ学習部410は、上記ステップ1011で設定し

10

20

30

40

50

た訓練データを深層ニューラルネットワークへ入力して訓練する（1012）。

【0201】

ステップ1011、1012の処理は、必ず1ステップ毎に行う必要はなく、数ステップ毎、またエピソード終了時といったタイミングで実施するように、条件分岐してもよい。

【0202】

次に、強化学習/エピソードデータ学習部410がエピソードの終了を判定する（1013）。強化学習/エピソードデータ学習部410は、エピソード終了フラグが設定されていればエピソードの終了値を判定する。判定がno（エピソード終了フラグの設定無し）の場合、ステップ1003に戻って上記処理を繰り返す。

【0203】

エピソード終了判定1013でyes（終了フラグの設定あり）の場合、ステップ1003からステップ1014のループを抜けて、ステップ1015へ進む。

【0204】

強化学習/エピソードデータ学習部410は、エピソード内のステップの繰り返しを終了したら、エピソードデータを出力してエピソードデータベース122に蓄積する（1015）。

【0205】

次に強化学習/エピソードデータ学習部410は、エピソードの連続成功回数を判定する（1016）。強化学習/エピソードデータ学習部410は、エピソードが1回成功して強化学習を終了するなら、エピソード成功判定はyesでよい。例えば10回連続成功を終了条件とした場合、10回に満たないなら判定結果はnoであるが、連続成功回数に1を加算してステップ1017に進む。エピソードが失敗ならば連続成功回数は0として、判定結果をnoとしてステップ1017に進む。

【0206】

次に強化学習/エピソードデータ学習部410は、実施したエピソードの回数を判定する（1017）。強化学習はエピソード毎の繰り返しであり、エピソードが成功回数判定（1016）でyesにならない限り、無限にステップ1001からステップ1018を繰り返すこととなる。そこで、強化学習/エピソードデータ学習部410は、エピソード回数の上限值を設け、エピソードの実施回数が上限値を超える場合にyesと判定してエピソード繰り返しのループを抜ける。この場合、強化学習は失敗となる。

【0207】

最後に強化学習の終了として、学習終了の処理を行う（1019）。強化学習/エピソードデータ学習部410は、強化学習過程のエピソードの情報などのログを出力する。

【0208】

以上が、図10のフローチャートによる、強化学習の一連の処理手順の一例の説明である。以上が、本発明における強化学習の内容についての説明である。

【0209】

<エピソード分析>

本発明における、エピソード分析の内容についての例を示す。

【0210】

エピソード分析は、図1、図3におけるエピソード分析部102及び図4のエピソード分析部430の内容に相当する。

【0211】

エピソード分析とは、エピソードデータベース419に蓄積されたエピソードデータを、強化学習過程におけるエピソードの成功、失敗の順を、他移動体の行動パターンについて個別に分析することで、学習パターンを構成する、という一連の処理である。

【0212】

個別のエピソードデータは図9A、図9Bに示した、1つのエピソードを構成するステップ901毎の状態量902と、行動903と、報酬904及び、状態905のデータである。1回の強化学習の結果として、強化学習過程について、1つ以上の、実際には多数

10

20

30

40

50

のエピソードデータがエピソードデータベース 419 に格納される。

【0213】

以降の説明におけるエピソードデータについては、艦船による追跡の強化学習についてのデータの例である。

【0214】

1つの強化学習過程についてのエピソードデータを図11A、図11Bに示す。図11Aは強化学習によってエピソードが成功した例を示し、図11Bは強化学習によってエピソードが失敗した例である。図示のグラフは横軸がエピソード(数)1101、縦軸がエピソード終了時のステップ数1102である。プロットの形は、凡例1103に示すように、見失いが黒四角(塗りつぶし四角)、拿捕が白三角(白抜き三角)、好位置が白丸(白抜き丸)である。また、シミュレーション結果が見失いと、拿捕の場合はエピソードの種別が失敗となり、シミュレーション結果が好位置の場合はエピソードの種別が失敗となる。

10

【0215】

図11Aのグラフでは、強化学習過程において150エピソード程度までは、ほぼ100ステップを超えずに見失いで、エピソードが失敗している。また、多少の拿捕のケースも見られる。150エピソード以降、エピソードは失敗であるが、ステップ数は増え始める。

【0216】

200エピソードを超えると、ステップ数が400ステップ達成の好位置でのエピソードの成功が見られるようになる。250エピソード以降はほぼエピソードの種別は成功であり、エピソードの成功が繰り返されて270エピソードで強化学習は成功として学習は完了した。

20

【0217】

失敗の例として図11Bのグラフでは、100エピソードまではほぼ100ステップを超えていないが、100エピソード以降、200エピソードまでで成功のエピソードが出現した。200エピソードから300エピソードの間は成功エピソードは多いが、シミュレーション結果が拿捕となった失敗エピソードも多い。

【0218】

300エピソード付近からは拿捕ばかりとなり、350エピソード以降、拿捕と見失いによる失敗エピソードばかりとなった。成功エピソードが連続することはなく、550エピソードで強化学習は失敗となった。

30

【0219】

他移動体の同一のシミュレーション結果で、エピソードが成功した強化学習の強化学習過程についての、エピソード失敗、成功別のステップ数のプロットを図12に示す。なお、以下の図では、上記図11Aと同様に、シミュレーション結果が見失いと、拿捕の場合はエピソードの種別が失敗となり、シミュレーション結果が好位置の場合はエピソードの種別が失敗となる。

【0220】

図12では、400ステップ達成のエピソード成功のプロットはエピソードの数が150から350、平均的には250エピソードで学習が完了している。

40

【0221】

見失い、又は拿捕で失敗したエピソードのステップ数も、強化学習過程のエピソード数に沿って徐々に多くなっている。但し100エピソードまでは150ステップ、平均的に見ればおよそ50ステップ程度までの増え方である。

【0222】

なお、プロットの状態を見ると強化学習の成功までの過程において、エピソードの種別のばらつきが大きい。強化学習では行動選択は乱択であり、深層ニューラルネットワーク10の訓練の結果にばらつきが大きく入るためと考えられる。理想的なエピソードを経た訓練ならば、エピソード数は少なくできると考えられる。

50

【0223】

学習パターンを検討するにあたり、見失いと、拿捕の失敗エピソードと400ステップ連続好位置の成功エピソードを、エピソード数とステップ数による単純な線図として、図13に示すように整理する。図中線1301は見失い、破線1302は拿捕を表す。これら2つの線は同じ傾向であるとした。網掛け線1303は成功エピソードである。図12における成功、失敗のプロットの平均的な値で描いた。

【0224】

強化学習を行った深層ニューラルネットワーク10がシミュレーションに成功する場合の、エピソード数の観点で理想的な、強化学習過程における失敗、成功エピソードの出現の範囲を図14に示す。図中、黒塗りつぶし1401は見失いを示し、斜め網掛け1402は拿捕を示し、点網掛け1403は連続400ステップで好位置を示す。

10

【0225】

エピソードデータ学習部420(122)は、150エピソードで強化学習が完了すると想定して、150エピソードまでの40エピソード分を成功のエピソードに割り当てる。強化学習過程の初期においては少ないステップ数での見失い、拿捕の失敗が多いので、ステップ数0から100の範囲で、20エピソードずつ見失い、拿捕の失敗のエピソードに割り当てる。継いでステップ数100から200の範囲で、20エピソードずつ見失い、拿捕の失敗を割り当てる。80エピソード目から成功のエピソードまでのエピソード数110の範囲を、ステップ数200から300、ステップ数300から400の見失い、拿捕のエピソードに割り当てる。

20

【0226】

なお、エピソードデータ学習部420(122)が、上記ステップ数の範囲とエピソード数の範囲に割り当てるシミュレーション結果の種別(見失い、拿捕、好位置)は、予め設定した比率(又は数)で割り当てることができる。

【0227】

図14に示した強化学習過程における成功、失敗のエピソードと、それらのステップ数の範囲より、図15に示す学習パターンが得られる。学習パターン1500は、エピソード結果1501と、終了ステップ数範囲の下限1502及び上限1503と、エピソード数1504で定義されるデータ構造である。エピソード結果1501はエピソード終了時の状態である。終了ステップ数範囲は下限1502と、上限1503で指定される。

30

【0228】

図15に示した学習パターン1500は、図16のグラフで示すように横軸にエピソード数、縦軸にステップ数としてプロットされる。図16では、ステップ数範囲は範囲の中心の値としている。

【0229】

以上により、他移動体の行動パターン(拿捕、見失い、好位置)に対して、エピソードデータベース419から収集したエピソードデータを用いて学習パターンを構成する例を示した。図4の行動制御計画装置105におけるエピソード分析部430は、エピソードデータ収集機能431、ステップ毎エピソード集計機能433、学習パターン構成機能434により処理される。

40

【0230】

エピソード分析は、図12の全てのエピソードデータの平均やばらつきなどを評価することにより、学習パターンの構成要素である、強化学習開始からのエピソード数に対するステップ数の範囲を決める処理であって、計算処理によって自動化できる。

【0231】

学習パターンを用いるエピソードデータ学習において、行動制御計画装置105はエピソードデータベース419からエピソードデータを収集する。以降でエピソードデータの収集方法について説明する。なお、エピソードデータの収集は他移動体のシミュレーション結果別に行う。

【0232】

50

この処理は、図3のフローチャートにおけるエピソード収集（学習用）311の内容であり、また図4の行動制御計画装置105の構成のエピソードデータ収集・一覧化機能422による処理である。

【0233】

図15に示した学習パターンに対応するエピソードデータのサンプリング結果を図17に示す。図17は、エピソードの結果順に、エピソード数の範囲毎にステップ数の範囲内のエピソードデータのサンプルが入っていることを示す。

【0234】

例えば、1エピソードから20エピソードまではステップ数100以下の見失い（図中四角）のサンプルを示し、21エピソードから40エピソードまではステップ数100以下の拿捕（図中三角）のサンプルを示し、また、111エピソードから140エピソードまではステップ数400以下の成功のエピソードのサンプル、となっている。

【0235】

図17の強化学習過程のエピソードの見失い、拿捕、好位置の分布は、図11Aに示した強化学習の成功例の分布とかなり異なっている。図17は見失い、拿捕、好位置のシミュレーション結果について、横軸（エピソード数）方向のエピソード範囲でのばらつきがないためである。

【0236】

そこで、図17のような学習パターンによるサンプリングにエピソード範囲でのばらつきを与えるために、強化学習過程におけるエピソード（の位置の）間の置換をする方法がある。エピソード間の置換とは、例えば、強化学習開始からの10エピソード目と20エピソード目のエピソードを交換する操作である。置換対象のエピソードの選択をランダムで決める置換をランダム置換と呼ぶ。

【0237】

図18Aに置換数50のランダム置換の例を示し、図18Bに置換数100のランダム置換の例を示し、強化学習過程についてのエピソード数とステップ数の分布を示す。図18A、図18Bでは、130エピソード以降のエピソードは置換の対象とはしていない。このように、学習パターンによるサンプリングにエピソード範囲でのばらつきを与えることができる。なお、置換という言葉で説明したが、シャッフルと言ってもよい。

【0238】

以上が学習パターンによるエピソードデータ収集方法についての説明である。

【0239】

以上が、本発明における、エピソード分析の内容についての例の説明である。

【0240】

<エピソードデータ学習>

本発明における、エピソードデータ学習の内容についての例を示す。

【0241】

エピソードデータ学習は、図1、図3におけるエピソードデータ学習部103及び図4のエピソードデータ学習部410の内容に相当する。

【0242】

エピソードデータ学習とは、他移動体のシミュレーション結果別の学習パターンに従ってエピソードデータを収集し、エピソード毎にエピソードのステップデータ（エピソードデータ）を直接に用いて深層ニューラルネットワークを訓練することで、状況に応じた行動の選択を学習する方法である。

【0243】

エピソードデータ学習にあたっては、まず条件を取得することは強化学習と同様である。但し、エピソードデータ学習のためには、他移動体のシミュレーション結果は複数であって、それらに対応する他移動体の条件、また学習パターンを取得することになる。

【0244】

エピソードデータ収集については、本実施例の学習パターンに対するエピソードデータ

10

20

30

40

50

の収集について説明した通りである。

【 0 2 4 5 】

エピソードデータ学習では、エピソードデータを直接に利用して、エピソードにおけるステップのデータを用いて深層ニューラルネットワークを訓練する。エピソードデータ学習の目的は、状態量に応じて行動を選択することであり、学習のモデルは、強化学習と同様にマルコフ決定過程 MDP のモデルに基づいている。

【 0 2 4 6 】

つまり上記強化学習の内容で示したように、方策関数と状態価値関数を深層ニューラルネットワークで定義し、エピソードの繰り返しで、各エピソードではステップ毎のデータを用いて深層ニューラルネットワークを訓練するものである。

10

【 0 2 4 7 】

深層ニューラルネットワークの構成、また訓練のための、例えば、上記(1)式から(10)式で表現される計算手法を用いる。また、ステップデータのサンプリングについても強化学習と同じである。

【 0 2 4 8 】

エピソードデータ学習と強化学習の違いは、強化学習はステップ毎にエージェントの行動選択に乱択と方策関数を利用した選択を混在させ、行動により状態を更新することにより状態量、行動、報酬のデータを得る。

【 0 2 4 9 】

これに対して、エピソードデータ学習は状態量、行動、報酬のデータを、例えばファイルから読み込むなど、直接に取得することである。エピソードデータ学習には、強化学習における探索と活用 (exploration and exploitation) といった特長はない。

20

【 0 2 5 0 】

この点でエピソードデータ学習そのものには、シミュレータや、場合によっては現実の環境における設備や実験は不要である。但し、学習が成功しない、無意味なエピソードデータは利用せず、事前に強化学習した際に得たエピソードデータを利用する。

【 0 2 5 1 】

ここで学習に必要なエピソードデータを学習パターンとして選別することで、学習を効率化し、さらに学習が難しくなる異なる環境や、他移動体のシミュレーション結果に対する学習を実現する手法である。

30

【 0 2 5 2 】

一連のエピソードデータ学習の処理手順の一例を、図19のフローチャートにより説明する。図19における処理の主体は行動制御計画装置105である。

【 0 2 5 3 】

まず、エピソードデータ学習部410は、エピソードデータをエピソードデータベース419(122)から収集する(1901)。

【 0 2 5 4 】

ステップ1902からステップ1918までの範囲はエピソード毎の繰り返しである。エピソードの繰り返しは学習パターンの指定に基づいて、エピソードデータベース122から収集されたエピソードデータを、強化学習過程における順序での繰り返しである。

40

【 0 2 5 5 】

なお、このエピソードデータには、他移動体の複数のシミュレーション結果のエピソードデータが混在している。例えば、シミュレーション結果が2つで、それぞれエピソードの数が同じである場合には、シミュレーション結果毎に一つずつのエピソードを並べる。1つのシミュレーション結果のエピソードの数が、もう一つの倍である場合、倍の方のエピソードの2つに対し、もう一つのエピソードを1つ並べるといった形での混在である。

【 0 2 5 6 】

エピソードの繰り返しにおいては、エピソードデータ学習部410は、まず環境、自移動体、他移動体の状態をリセットする(1903)。

50

【0257】

ステップ1904からステップ1909の範囲は1エピソードにおけるステップ毎の繰り返しである。

【0258】

エピソードデータ学習部410は、現在のステップの状態量、行動、報酬を取得する(1905)。強化学習では、図10に示したように、ステップの繰り返しでは、行動選択、行動実行、状態に基づく失敗、成功の判定があったが、エピソードデータ学習では、ステップのデータを読み込むだけである。

【0259】

次にエピソードデータ学習部410は、DNN訓練機能414を利用してステップのデータで深層ニューラルネットワークを訓練する(1906)。

10

【0260】

エピソードデータ学習部410は、エピソードの終了を判定する(1907)。エピソードデータでステップ数が決まっているので最後のステップであればエピソードの終了となる。一方、終了でなければステップを更新し(1908)、次のステップで上記処理を繰り返す。

【0261】

エピソードデータ学習部410は、エピソードの処理が終了する度に、学習の完了のテストをステップ1910~1915で行う。

【0262】

ステップ1910からステップ1915の範囲は、全てのシミュレーション結果についての繰り返しである。

20

【0263】

そして、ステップ1911からステップ1914の範囲は、指定された試行回数についての繰り返しである。この試行回数が、例えば5回である場合、ある一つのシミュレーション結果について5回の試行を行うということであり、連続して5回成功ならば、そのシミュレーション結果に対する学習は成功したということの意味する。

【0264】

まず、エピソードデータ学習部410は、訓練結果(学習結果)の深層ニューラルネットワークを用いて指定されたエピソードを実行する(1912)。これは図3におけるステップ305と同様の処理であって、初期状態から1ステップ毎に方策関数により行動を選択し、状態を更新することである。

30

【0265】

次に、エピソードデータ学習部410は、エピソードの失敗か成功かを判定する(1913)。失敗した場合がyesであって、直ちにステップ1910からステップ1915の繰り返しを抜ける。失敗しない場合(no)には、繰り返しを継続する。全ての繰り返しで成功したなら、ステップ1913からステップ1910の繰り返しを抜けることはない。

【0266】

次に、エピソードデータ学習部410は、全てのシミュレーション結果で成功したかを判定する(1916)。成功した場合がyesであり、この場合には学習が成功でありエピソードの繰り返しから抜ける。

40

【0267】

次にエピソード回数を判定する(1917)。この処理は、所定のエピソード繰り返し回数の間にエピソードが成功しない場合の判定であり、判定がyesの場合、学習失敗としてエピソードの繰り返しから抜ける。全てのエピソードの繰り返しで学習を行う場合には、ステップ1917の判定は不要である。

【0268】

最後にエピソードデータ学習部410は、学習終了の処理を行う(1919)。この処理は、ログなどの出力処理である。

50

【0269】

以上が、図19のフローチャートによる、一連のエピソードデータ学習の処理手順の例の説明である。上記処理によって、他移動体の複数のシミュレーション結果（例えば、拿捕、見失い、好位置など）について予め設定された学習パターンに一致するエピソードデータを抽出してエピソードデータ学習を行って、学習結果をエピソードデータベース122へ蓄積することができる。

【0270】

これにより、強化学習で算出された全てのエピソードデータベース122を用いるのではなく、予め設定された学習パターンで抽出したエピソードデータをエピソードデータベース122から抽出することで、複数のシミュレーション結果のエピソードデータで深層ニューラルネットワーク10を直接訓練する。

10

【0271】

エピソードデータ学習では学習パターンの設定に応じて、効率的に深層ニューラルネットワークの訓練が可能となり、例えば、少ないエピソード数で訓練を完了させながらも、学習が成功しやすくなる。これにより、深層ニューラルネットワーク10の行動制御の精度を向上させて、自律的に行動する移動体の運用を容易にすることが可能となる。

【0272】

以上が、本発明におけるエピソードデータ学習の内容についての例の説明である。

【0273】

本発明の移動体の行動制御計画装置は、計算処理を主な技術としており、ソフトウェアを主な実現手段とする。したがってパッケージソフト、クラウドサービスといった形態での実現が考えられる。学習した結果の方策関数、状態価値関数は制御ソフトウェアとして機械に組み込める。よって、実際に行動する艦艇、車両といった移動体に組み込まれたり、また移動体を遠隔的に制御する他の移動体や装置、施設での活用も考えられる。本発明はそのような活用の範囲を限定するものではない。

20

【0274】

以上が、他移動体の様々な行動に対する自移動体の行動制御のために、他移動体の行動に対して強化学習を行い、その結果の様々な行動のエピソードデータを用いたエピソードデータ学習によって、自移動体の行動制御を計画する技術の説明である。

【0275】

<結び>

以上説明したように、本発明の行動制御計画装置105は、自移動体に対して他移動体の行動が多様な関係となる場合でも、例えば追跡、検知といった目的を達成するための、自移動体の行動制御を深層ニューラルネットワークで学習し、学習結果の深層ニューラルネットワークを自律的に行動可能な移動体の制御に適用することが可能となる。

30

【0276】

また、エピソード中のステップのデータを、深層ニューラルネットワークの訓練におけるパラメータの最適化のため、目的関数である損失関数値に基づきエピソードデータに優先付けすることで、効率的に強化学習でき、強化学習を行った深層ニューラルネットワークではエピソードが成功する可能性を向上させることが可能となる。

40

【0277】

個別の環境、特に他移動体のシミュレーション結果に対しては強化学習を実施した深層ニューラルネットワーク10がシミュレーションで成功する。成功した強化学習のエピソードをエピソードデータとして蓄積する。複数の他移動体のシミュレーション結果のエピソードデータを用いて、エピソードのステップデータで直接に深層ニューラルネットワークの訓練を行うエピソードデータ学習により、複数のシミュレーション結果に対応する行動制御を深層ニューラルネットワークに学習させることができる。

【0278】

エピソードデータ学習の成功のために、蓄積されたエピソードデータを収集して分析する。個別の環境、他移動体のシミュレーション結果の種類別に、成功した強化学習につい

50

て、エピソードの順序に関して成功と失敗のエピソードとそのステップ数を一覧化することで、エピソード順序に対して失敗、成功とステップ数の範囲を学習パターンとして定義する。

【0279】

上記定義されたエピソード順序に対するエピソードデータをエピソードデータ学習に利用することでエピソードデータ学習を実施することで、自移動体に対して他移動体の行動が多様な関係となる場合でもエピソードが成功する可能性を向上させることができる。

【0280】

強化学習、エピソードデータ学習により移動体の行動制御を計画すれば、マニュアルによる計画が困難な多様な環境、他移動体の行動パターンに対応することが可能となる。さらには、本実施例の行動制御計画装置は、人間が考え付くことが困難と考えられる行動制御内容が得られる可能性がある。

10

【0281】

なお、上記実施例において、行動制御計画装置105で学習した深層ニューラルネットワークを適用する移動体として哨戒又は追跡を行う艦船の行動制御に本発明を適用する例を示したが、これに限定されるものではなく、航空機や車両、潜水艦或いはドローン等で自律的に行動する移動体に適用することができる。

【0282】

また、上記実施例において、行動制御計画装置105は、シミュレータ400（又は112）を含む例を示したが、これに限定されるものではなく、シミュレータ400を外部の計算機で実施してもよい。この場合、行動制御計画装置105は、図1の条件設定111で受け付けたシミュレーション条件を、外部のシミュレータに送信してシミュレーションを実行させ、シミュレーションの結果を取得する。この例では、行動制御計画装置105と外部シミュレータで構成された行動制御計画システムとすることができる。そして、行動制御計画装置105は、外部からのシミュレーション結果に基づいて深層ニューラルネットワーク10の強化学習を実施する。

20

【0283】

以上のように、上記実施例の行動制御計画装置105は、以下のような構成とすることができる。

【0284】

(1) プロセッサ(11)とメモリ(12)を有して、他移動体の行動に対して自移動体の行動を計画する行動制御計画装置(105)であって、前記自移動体と他移動体の行動をシミュレーション条件として設定するシミュレーション条件設定部(111)と、前記シミュレーション条件と予め設定された機械学習モデル(深層ニューラルネットワーク10)に基づいて、前記自移動体(504)と他移動体(505)の位置関係を所定の時間間隔毎にステップデータ(900A)として出力するシミュレータ(112)と、前記ステップデータ(900A)を取得して、前記機械学習モデル(10)に前記ステップデータ(900A)を与えて学習結果として前記自移動体と他移動体の前記位置関係を示す状態量(902)と報酬(904)と状態(905)を含むエピソードデータを出力させ、当該エピソードデータをエピソードデータ蓄積部(エピソードデータベース122)に蓄積する強化学習部(101)と、前記エピソードデータ蓄積部(122)に蓄積された前記エピソードデータを分析して学習パターンを生成するエピソード分析部(102)と、前記学習パターンに該当する前記エピソードデータを前記エピソードデータ蓄積部(122)から取得して学習用エピソードデータを生成(エピソードデータ収集・一覧化機能422)し、前記機械学習モデル(10)に前記学習用エピソードデータを与えて学習させるエピソードデータ学習部(103)と、を有することを特徴とする行動制御計画装置。

30

40

【0285】

上記構成により、行動制御計画装置105は、様々な環境に対して、特に自移動体に対して他移動体の行動が多様な関係となる場合でも、例えば追跡、検知といった目的を達成するための、自移動体の行動を学習できる。

50

【0286】

(2) 上記(1)に記載の行動制御計画装置(105)であって、前記強化学習部(101)は、前記ステップデータ(900A)の開始から所定の終了条件に達するまでのステップデータ(900A)を1つのエピソードとし、前記エピソードは前記ステップデータ(900A)に対応するステップを前記エピソードデータに設定し、前記エピソードデータの前記状態(905)には前記ステップデータ(900A)毎のシミュレーション結果を設定し、前記エピソードには前記シミュレーション結果に応じて所定の目的を達成したか否かを示すエピソード種別(1007、1009)を設定し、前記エピソード分析部(102)は、前記エピソード種別と前記シミュレーション結果と前記エピソードの数と前記エピソードデータのステップ数に基づいて前記学習パターンを生成することを特徴とする行動制御計画装置。

10

【0287】

上記構成により、エピソード分析部102は、全てのエピソードデータの平均やばらつきなどを評価することにより、学習パターンの構成要素である、強化学習開始からのエピソード数に対するステップ数の範囲を絞り込むことで、学習対象の数を抑制しながら効率よく機械学習モデルの学習(訓練)を実現することが可能となる。

【0288】

(3) 上記(2)に記載の行動制御計画装置(105)であって、前記エピソード分析部(102)は、前記シミュレーション結果毎に前記エピソードの数の範囲(1504)と、前記エピソードデータのステップ数の範囲(1502、1503)を設定し、前記エピソードデータ学習部(103)は、前記エピソードデータ蓄積部(122)から、前記シミュレーション結果毎に前記エピソードの数の範囲(1504)と、前記エピソードデータのステップ数の範囲(1502、1503)に該当するエピソードを取得して前記機械学習モデル(10)に学習させることを特徴とする行動制御計画装置。

20

【0289】

上記構成により、エピソードデータ学習部103は、学習パターンの設定に応じて、効率的に深層ニューラルネットワーク10の訓練が可能となり、例えば、少ないエピソード数で訓練を完了させながらも、学習が成功しやすくなる。これにより、深層ニューラルネットワーク10の行動制御の精度を向上させて、自律的に行動する移動体の運用を容易にすることが可能となる。

30

【0290】

(4) 上記(3)に記載の行動制御計画装置(105)であって、前記エピソードデータ学習部(103)は、前記シミュレーション結果毎に、前記エピソードの数(1504)とステップ数の範囲(1502、1503)で抽出したエピソードを、前記エピソードの順序で所定の置換を行うことを特徴とする行動制御計画装置。

上記構成により、エピソード分析部102は、エピソードの順序で所定の置換を行うことで、深層ニューラルネットワーク10の訓練結果のばらつきを抑制することができる。

【0291】

(5) 上記(3)に記載の行動制御計画装置(105)であって、前記エピソードデータ学習部(103)は、前記取得したエピソードに対応する前記エピソードデータを前記機械学習モデル(10)に直接入力して学習させる(1905)ことを特徴とする行動制御計画装置。

40

【0292】

上記構成により、エピソードデータ学習部103は、エピソードデータを直接用いて深層ニューラルネットワーク10を訓練することで、状況に応じた行動の選択を学習することが可能となる。

【0293】

(6) 上記(1)に記載の行動制御計画装置(105)であって、前記所定の目的は、前記自移動体(504)が前記他移動体(505)の追跡であることを特徴とする行動制御計画装置。

50

【0294】

上記構成により、自移動体504が他移動体505を見失うことなく、かつ、他移動体505に拿捕されることなく追跡を行う深層ニューラルネットワーク10を生成することができる。

【0295】

(7)上記(2)に記載の行動制御計画装置(105)であって、前記エピソード分析部(102)は、前記エピソードの数の範囲(1504)と、前記エピソードデータのステップ数の範囲(1502、1503)において、前記シミュレーション結果の種別に応じた所定の比率で前記エピソードを割り当てることを特徴とする行動制御計画装置。

【0296】

エピソード分析部102は、エピソード数の範囲と、エピソードデータのステップ数の範囲内で、シミュレーション結果の種別毎に割り当てる比率を設定することで、エピソード数とステップ数に応じたシミュレーション結果の種別の傾向を適用することが可能となる。

【0297】

なお、本発明は上記した実施例に限定されるものではなく、様々な変形例が含まれる。例えば、上記した実施例は本発明を分かりやすく説明するために詳細に記載したものであり、必ずしも説明した全ての構成を備えるものに限定されるものではない。また、ある実施例の構成の一部を他の実施例の構成に置き換えることが可能であり、また、ある実施例の構成に他の実施例の構成を加えることも可能である。また、各実施例の構成の一部について、他の構成の追加、削除、又は置換のいずれもが、単独で、又は組み合わせても適用可能である。

【0298】

また、上記の各構成、機能、処理部、及び処理手段等は、それらの一部又は全部を、例えば集積回路で設計する等によりハードウェアで実現してもよい。また、上記の各構成、及び機能等は、プロセッサがそれぞれの機能を実現するプログラムを解釈し、実行することによりソフトウェアで実現してもよい。各機能を実現するプログラム、テーブル、ファイル等の情報は、メモリや、ハードディスク、SSD(Solid State Drive)等の記録装置、又は、ICカード、SDカード、DVD等の記録媒体に置くことができる。

【0299】

また、制御線や情報線は説明上必要と考えられるものを示しており、製品上必ずしも全ての制御線や情報線を示しているとは限らない。実際には殆ど全ての構成が相互に接続されていると考えてもよい。

【符号の説明】

【0300】

- 11 プロセッサ
- 12 メモリ
- 13 ストレージ装置
- 101 強化学習部
- 102 エピソード分析部
- 103 エピソードデータ学習部
- 104 オペレータ
- 105 行動制御計画装置
- 111 条件設定
- 112 シミュレータ
- 113 状態更新
- 114 行動選択
- 115 学習器
- 116 訓練データ取得

10

20

30

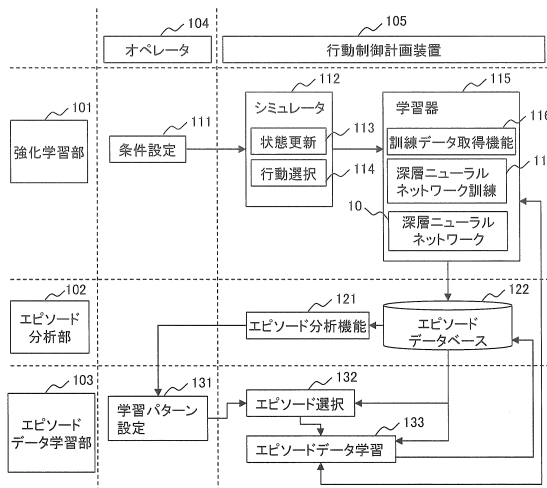
40

50

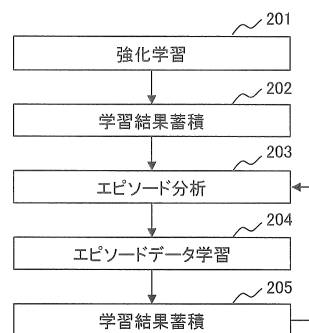
- 1 1 7 深層ニューラルネットワーク訓練
- 1 2 1 エピソード分析
- 1 2 2 エピソードデータベース
- 1 3 1 学習パターン設定
- 1 3 2 エピソード選択
- 1 3 3 エピソードデータ学習
- 4 0 0 シミュレータ
- 4 0 1 行動環境設定機能
- 4 0 2 移動体行動処理機能
- 4 0 3 移動体状態量算出機能 10
- 4 0 4 移動体状態判定機能
- 4 1 0 強化学習 / エピソードデータ学習部
- 4 1 1 方策・状態価値関数構成機能
- 4 1 2 エピソード処理機能
- 4 1 3 報酬設定機能
- 4 1 4 DNN訓練機能
- 4 1 5 DNN訓練用データサンプリング機能
- 4 1 6 方策・状態価値算出機能
- 4 1 7 エピソードデータ蓄積機能
- 4 1 8 行動環境データベース 20
- 4 1 9 エピソードデータベース
- 4 2 0 エピソードデータ学習部
- 4 2 1 学習パターン取得機能
- 4 2 2 エピソードデータ収集・一覧化機能
- 4 2 3 エピソードデータ学習機能
- 4 3 0 エピソード分析部
- 4 3 1 エピソードデータ収集機能
- 4 3 2 エピソードデータ描画機能
- 4 3 3 ステップ毎エピソード集計機能
- 4 3 4 学習パターン構成機能 30

【図面】

【図 1】

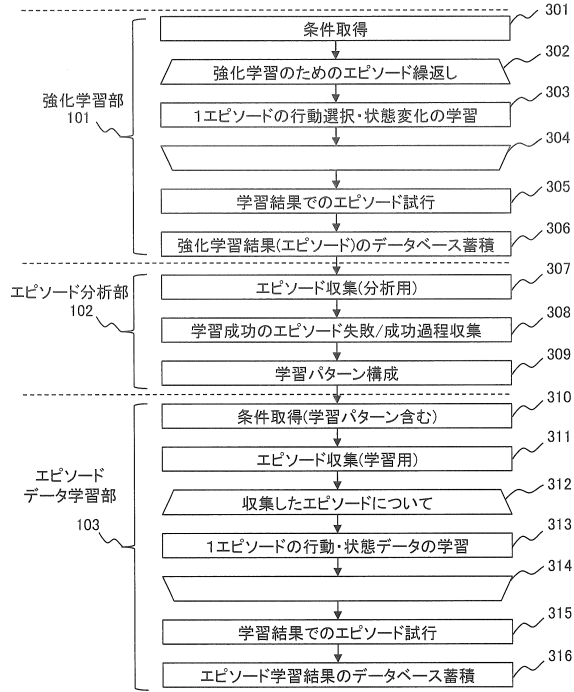


【図 2】

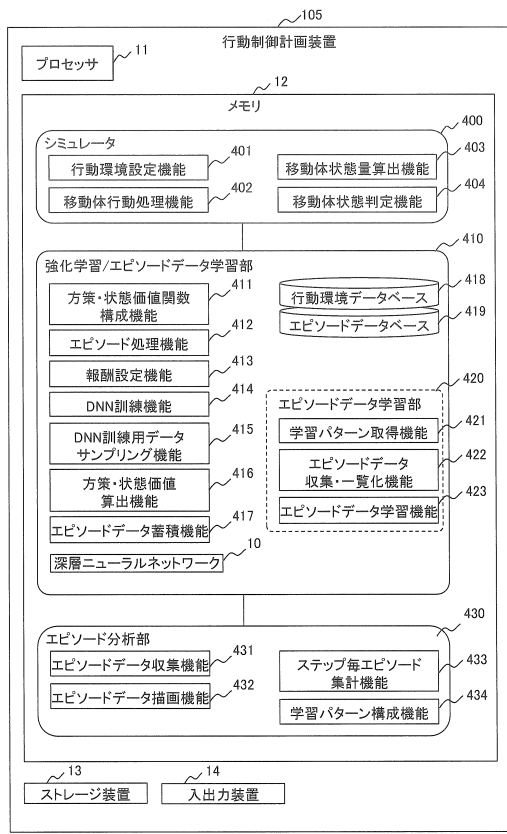


40

【図3】



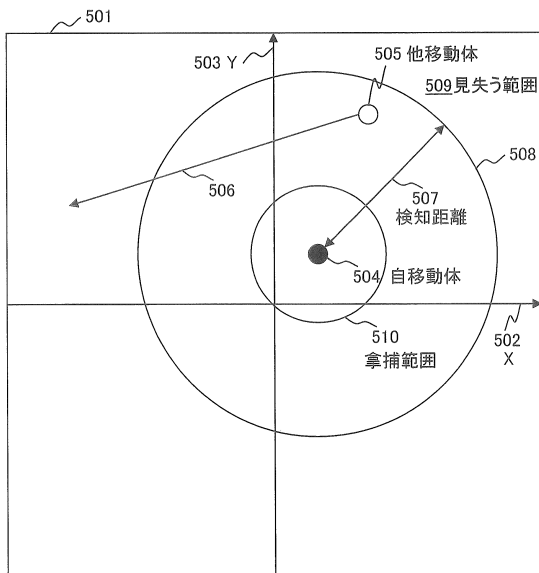
【図4】



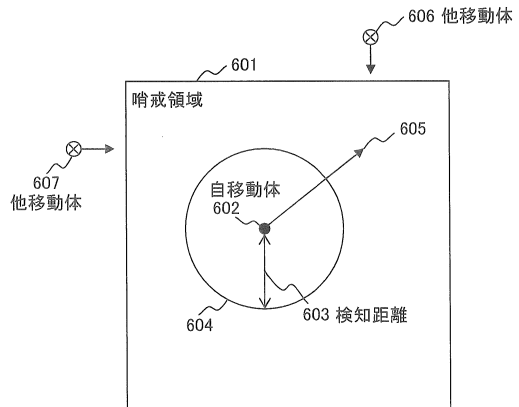
10

20

【図5】



【図6】

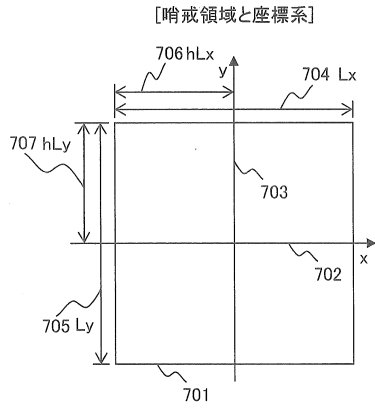


30

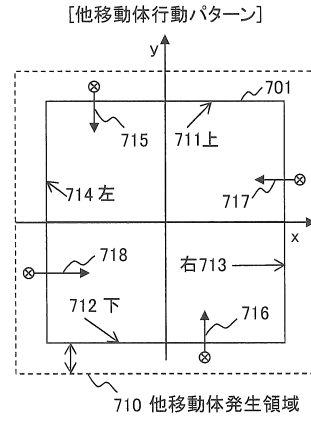
40

50

【図7A】

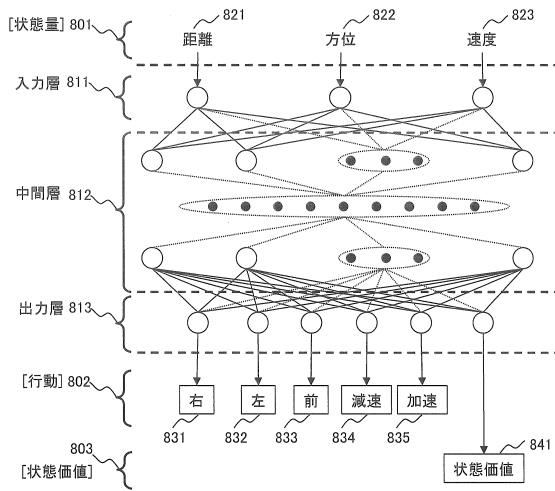


【図7B】



10

【図8】



【図9A】

900A ステップデータ(失敗例)

ステップ	状態量			行動					報酬	状態
	距離	方位	速度	前	右	左	減速	加速		
1	0.2	0.7	0.3	1	0	0	0	0	0	好位置
2	0.23	0.68	0.3	0	1	0	0	0	0	好位置
...										
273	0.5	0.12	0.1	1	0	0	0	0	-1	見失い

911 912 913 914 915 916 917 918

20

30

40

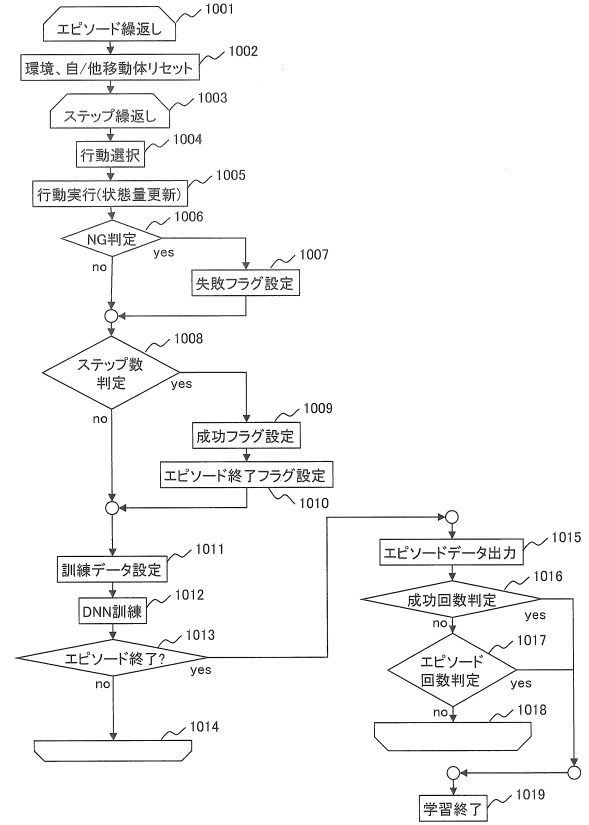
50

【図 9 B】

900B ステップデータ(成功例)

ステップ	状態量			行動					報酬	状態
	距離	方位	速度	前	右	左	減速	加速		
1	0.2	0.7	0.3	1	0	0	0	0	0	好位置
2	0.23	0.68	0.3	0	1	0	0	0	0	好位置
...										
400	0.18	0.08	0.4	0	0	0	1	0	1	好位置

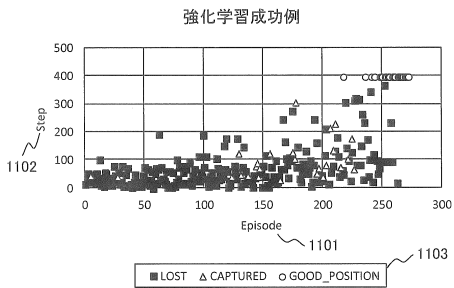
【図 1 0】



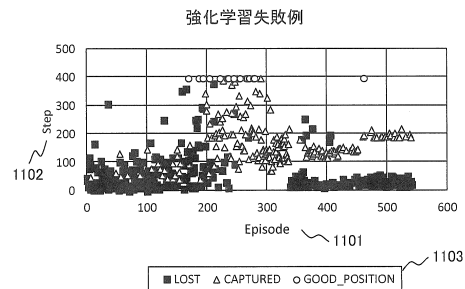
10

20

【図 1 1 A】



【図 1 1 B】

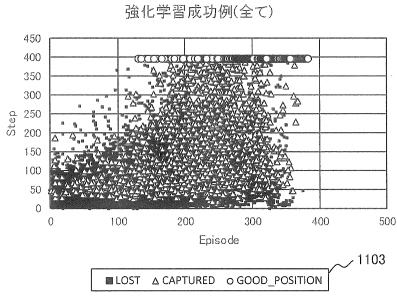


30

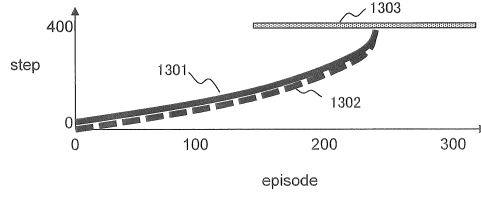
40

50

【図 1 2】

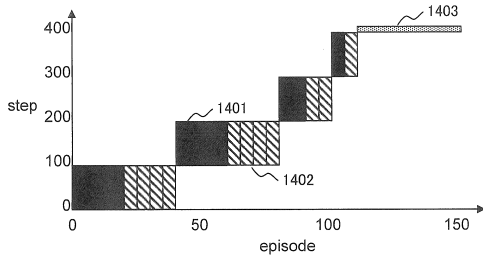


【図 1 3】



10

【図 1 4】



【図 1 5】

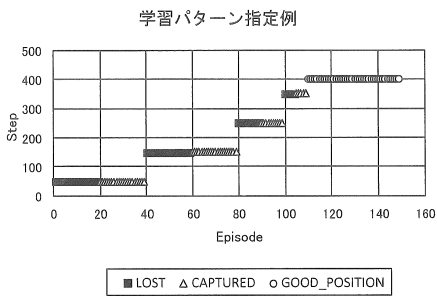
1500 学習パターン

エピソード結果	終了ステップ数範囲		エピソード数
	下限	上限	
見失い	1	100	20
拿捕	1	100	20
見失い	101	200	20
拿捕	101	200	20
見失い	201	300	10
拿捕	201	300	10
見失い	301	400	5
拿捕	301	400	5
好位置	400	400	40

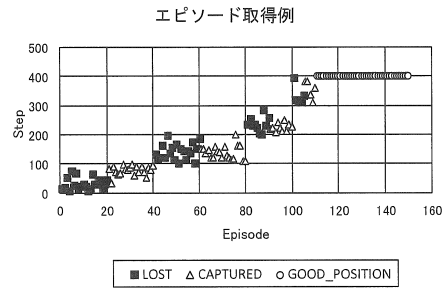
Labels 1501, 1502, 1503, 1504 are positioned below the table.

20

【図 1 6】



【図 1 7】

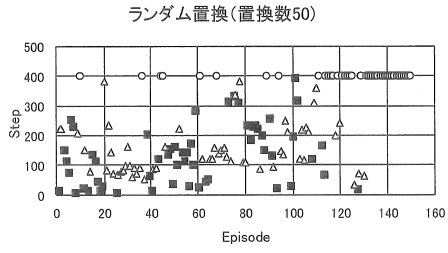


30

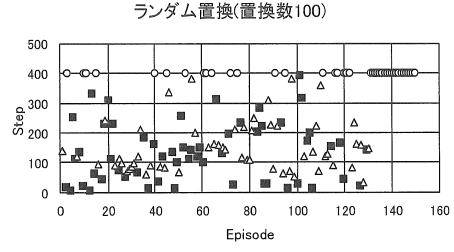
40

50

【図18A】

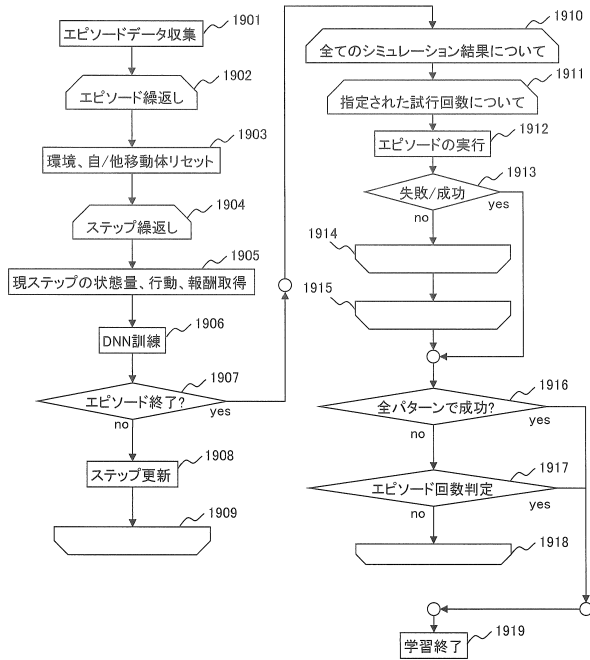


【図18B】



10

【図19】



20

30

40

50

フロントページの続き

- (56)参考文献 特開2021-18484(JP,A)
特開2020-190854(JP,A)
特開2020-162438(JP,A)
特表2019-529135(JP,A)
韓国登録特許第10-2169876(KR,B1)
松本耕平,外3名,“予測状態表現に基づく深層強化学習を用いた動的環境下における移動ロボットナビゲーション”,第38回日本ロボット学会学術講演会予稿集DVD-ROM 2020年,日本,一般社団法人日本ロボット学会,2020年10月09日,p.1-4
- (58)調査した分野 (Int.Cl.,DB名)
G06N 20/00