

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2019-82809

(P2019-82809A)

(43) 公開日 令和1年5月30日(2019.5.30)

(51) Int. Cl.	F I	テーマコード (参考)
G08G 1/00 (2006.01)	G08G 1/00 C	5H181
G08G 1/09 (2006.01)	G08G 1/09 F	
G06N 20/00 (2019.01)	G06N 99/00 150	

審査請求 未請求 請求項の数 8 O L (全 17 頁)

(21) 出願番号	特願2017-209276 (P2017-209276)	(71) 出願人	000004226 日本電信電話株式会社 東京都千代田区大手町一丁目5番1号
(22) 出願日	平成29年10月30日(2017.10.30)	(74) 代理人	110001519 特許業務法人太陽国際特許事務所
		(72) 発明者	幸島 匡宏 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
		(72) 発明者	堤田 恭太 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内
		(72) 発明者	松林 達史 東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内

最終頁に続く

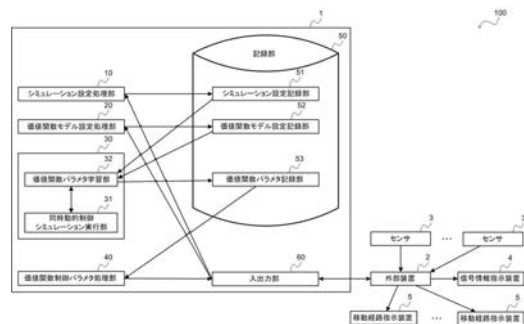
(54) 【発明の名称】 価値関数パラメタ学習装置、信号情報指示装置、移動経路指示装置、価値関数パラメタ学習方法、信号情報指示方法、移動経路指示方法、およびプログラム

(57) 【要約】

【課題】 移動体の数が増加しても、最適な交通状況を実現することができるようにする。

【解決手段】 同時動的制御シミュレーション実行部31が、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行し、価値関数パラメタ学習部32が、同時動的制御シミュレーション実行部31によるシミュレーションの結果に基づいて、価値関数のパラメタを学習する。

【選択図】 図1



【特許請求の範囲】**【請求項 1】**

交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行する同時動的制御シミュレーション実行部と、

前記同時動的制御シミュレーション実行部によるシミュレーションの結果に基づいて、前記価値関数のパラメタを学習する価値関数パラメタ学習部と、

を含む価値関数パラメタ学習装置。

10

【請求項 2】

請求項 1 記載の価値関数パラメタ学習装置によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各信号機に対して指示を行う信号情報指示装置

。

【請求項 3】

請求項 1 記載の価値関数パラメタ学習装置によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各区間を通過する移動体に対して行う進むべき経路の指示を行う移動経路指示装置。

20

【請求項 4】

同時動的制御シミュレーション実行部が、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行するステップと、

価値関数パラメタ学習部が、前記同時動的制御シミュレーション実行部によるシミュレーションの結果に基づいて、前記価値関数のパラメタを学習するステップと、

を含む価値関数パラメタ学習方法。

【請求項 5】

信号情報指示装置が、請求項 4 記載の価値関数パラメタ学習方法によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各信号機に対して指示を行う信号情報指示方法。

30

【請求項 6】

移動経路指示装置が、請求項 4 記載の価値関数パラメタ学習方法によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各区間を通過する移動体に対して行う進むべき経路の指示を行う移動経路指示方法。

【請求項 7】

コンピュータを、請求項 1 記載の価値関数パラメタ学習装置の各部として機能させるためのプログラム。

40

【請求項 8】

コンピュータを、請求項 2 記載の信号情報指示装置として機能させるためのプログラム

。

【発明の詳細な説明】**【技術分野】****【0001】**

本発明は、価値関数パラメタ学習装置、信号情報指示装置、移動経路指示装置、価値関数パラメタ学習方法、信号情報指示方法、移動経路指示方法、およびプログラムに関し、

50

特に、移動体による交通を制御するための価値関数パラメタ学習装置、信号情報指示装置、移動経路指示装置、価値関数パラメタ学習方法、信号情報指示方法、移動経路指示方法、およびプログラムに関する。

【背景技術】

【0002】

従来から、都市の交通渋滞や大規模イベントなどにおける人の混雑は社会的な課題になっている。交通渋滞は、渋滞中の車に乗車する人の時間を奪い、流通システムの遅れを生む原因にもなる。イベント会場における混雑もドミノ倒しなどの非劇的な雑踏事故を生む原因になりうる。

【0003】

この点、強化学習（非特許文献1）によって信号機の制御を行うことで、車両の待ち時間を減少させる技術が存在する（非特許文献2）。

【先行技術文献】

【非特許文献】

【0004】

【非特許文献1】Reinforcement learning: An introduction, Richard S Sutton and Andrew G. Barto, MIT press Cambridge, 1998.

【非特許文献2】Using a deep reinforcement learning agent for traffic signal control, Genders, Wade and Razavi, Saiedeh, arXiv preprint arXiv:1611.01142, 2016.

【発明の概要】

【発明が解決しようとする課題】

【0005】

実際には目的地に到達するための移動経路は多数存在しているため、車両や人の移動経路についても加味して考える必要がある。

【0006】

しかし、非特許文献2の技術では、移動経路について考慮されていないため、移動経路を加味した渋滞緩和を行うことができなかった。

【0007】

また、移動経路について強化学習を行ったとしても、各車両毎に移動経路を決定するため、車両数が増大すると、アクションの数が指数的に増大してしまう。このような膨大なアクション数を持つ場合、探索空間が増大し、強化学習によって正しく各車両の移動経路を推定することが極めて困難となる、という問題があった。

【0008】

本発明は、上記の点に鑑みてなされたものであり、移動体の数が増加しても、最適な交通状況を実現するための価値関数パラメタを学習することができる価値関数パラメタ学習装置、価値関数パラメタ学習方法、およびプログラムを提供することを目的とする。

【0009】

また、本発明は、移動体の数が増加しても、最適な交通状況を実現することができる信号情報指示装置、移動経路指示装置、信号情報指示方法、移動経路指示方法、およびプログラムを提供することを目的とする。

【課題を解決するための手段】

【0010】

本発明に係る価値関数パラメタ学習装置は、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行する同時動的制御シミュレーション実行部と、前記同時動的制御シミュレーション実行部によるシミュレーションの結果に基づいて、前記価値関数のパラメタを学習する価値関数パラメタ学習部とを備えて構成される。

【0011】

10

20

30

40

50

また、本発明に係る価値関数パラメタ学習方法は、同時動的制御シミュレーション実行部が、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行するステップと、価値関数パラメタ学習部が、前記同時動的制御シミュレーション実行部によるシミュレーションの結果に基づいて、前記価値関数のパラメタを学習するステップとを含む。

【0012】

本発明に係る価値関数パラメタ学習装置及び価値関数パラメタ学習方法によれば、同時動的制御シミュレーション実行部が、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行する。

10

【0013】

そして、価値関数パラメタ学習部が、同時動的制御シミュレーション実行部によるシミュレーションの結果に基づいて、価値関数のパラメタを学習する。

【0014】

このように、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを行い、シミュレーションの結果に基づいて、価値関数パラメタを学習することにより、移動体の数が増加しても、最適な交通状況を実現するための価値関数パラメタを学習することができる。

20

【0015】

本発明に係る信号情報指示装置は、価値関数パラメタ学習装置によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各信号機に対して指示を行うことができる。

【0016】

また、本発明に係る信号情報指示方法は、信号情報指示装置が、上記価値関数パラメタ学習方法によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各信号機に対して指示を行うことができる。

30

【0017】

本発明に係る移動経路指示装置は、上記価値関数パラメタ学習装置によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各区間を通過する移動体に対して行う進むべき経路の指示を行うことができる。

【0018】

また、本発明に係る移動経路指示方法は、移動経路指示装置が、上記価値関数パラメタ学習方法によって学習された前記価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する前記状態について、前記行動を決定し、前記決定された行動に従って、各区間を通過する移動体に対して行う進むべき経路の指示を行うことができる。

40

【0019】

本発明に係るプログラムは、上記の価値関数パラメタ学習装置の各部として機能させるためのプログラムである。

【0020】

また、本発明に係るプログラムは、上記の信号情報指示装置の各部として機能させるためのプログラムである。

50

【発明の効果】

【0021】

本発明の価値関数パラメタ学習装置、価値関数パラメタ学習方法、およびプログラムによれば、移動体の数が増加しても、最適な交通状況を実現するための価値関数パラメタを学習することができる。

【0022】

本発明の信号情報指示装置、移動経路指示装置、信号情報指示方法、移動経路指示方法、およびプログラムによれば、移動体の数が増加しても、最適な交通状況を実現することができる。

【図面の簡単な説明】

10

【0023】

【図1】本発明の実施の形態に係る交通制御システムの構成を示す概略図である。

【図2】本発明の実施の形態に係る価値関数パラメタ学習処理ルーチンを示すフローチャートである。

【図3】本発明の実施の形態に係る制御処理ルーチンを示すフローチャートである。

【図4】本発明の実施の形態に係る提案エージェントの一例を表す図である。

【発明を実施するための形態】

【0024】

以下、本発明の実施の形態について図面を用いて説明する。

【0025】

20

<本発明の実施の形態の原理>

まず、本発明の実施形態の原理について説明する。

【0026】

本実施形態は、都市の交通渋滞や大規模イベントなどにおける人の混雑を解消するための技術である。特に本発明の実施形態は人や車などの「移動体」の経路とそれら移動体に対して進行許可・停止などの指示を与える「信号機」を同時に最適化することで混雑を解消する技術である。本発明の実施形態の新規性は、移動体の経路と信号機の最適化を個別に行うのではなく、同時に行うことで全体最適化を行うことにある。

【0027】

この全体最適のための本実施形態の鍵となるアイデアは、移動体の経路と信号機を制御する指示主体（これをエージェント（Agent）と呼ぶ）の定義とこの指示主体の最適な制御策の推定方法にある。本実施形態に係るエージェント（以下、提案エージェント）の定義と学習方法によって、移動体の数や指示決定に用いる観測のセンサ数が膨大となる場合であっても、最適な制御策の推定が可能になる。これによって、人の混雑を解消することを実現する。

30

【0028】

本実施形態は、このような車、人などの混雑を解消する混雑緩和を行うための技術である。特に本実施形態は人や車などの「移動体の経路」とそれら移動体に対して進行許可・停止などの指示を与える「信号機」を同時に最適化することで混雑を解消する技術である。

40

【0029】

なお、ここで移動体と呼んでいるものは、人や動物、バイク、車、鉄道、ヘリコプター、飛行機など、一般に移動する、動くもの全てのものを指す。川や水路を流れる水や、ネットワークを流れるパケットも移動体である。

【0030】

また、信号機と呼んでいるものは、いわゆる道路の信号機（交通信号機）だけでなく、警察官による手信号など、人、車などの上記移動体に対して進行許可・停止などを指示する機能をもつ全てのものを指す。

【0031】

以下、これより簡便さのために「移動体の経路」とは車の移動経路、「信号機」とは道

50

路の交通信号機を意味するものとして交通制御の文脈で記述を進めるが、「移動体の経路」と「信号機」は上記のいずれのものであってもよい。

【0032】

本実施形態では、道路に設置してあるセンサや車、信号から観測される値に基づいて、動的（適応的）に信号機と車両の移動経路を制御できる、という設定を考える。

【0033】

ここで、制御とは、信号機の場合には信号を青にするもしくは赤にすること、経路の場合には車両が移動する経路を強制し、もしくは間接的に指示することを指す。なお、センサなどの観測値に基づいて信号を変えることのできる信号は感応式信号などと呼ばれ普及している。

【0034】

上記の設定のもと、本実施形態は、任意の観測値を入力、とるべき制御策を出力とする関数（この関数のことを最適方策と呼ぶ）をシミュレーションまたは実環境を通して学習する。一旦最適方策が学習できれば、それに従って制御を実施することで渋滞が緩和できる。

【0035】

このような最適方策を学習するアプローチは強化学習と呼ばれる。強化学習は、本発明のような信号機と車両の移動経路のように複数のものを制御する場合ではなく、信号機単独のものを制御する場合でも利用できるアプローチである。実際、強化学習によって信号機単独の制御を行う既存技術が存在する（非特許文献2）。

【0036】

本発明の実施形態はその既存技術の大きな発展系の一つである。そこで、まずこの既存研究について紹介する。

【0037】

<<強化学習>>

まずはじめに強化学習について簡単に説明する。強化学習はマルコフ決定過程（Markov Decision Process、MDP）（非特許文献1）として定義された設定で最適方策を見つける手法である。MDPは、簡単にいえば行動主体（例えばロボット）と外界の相互作用を記述したものであり、ロボットがとりうる状態の集合

$$\mathcal{S} = \{s_1, s_2, \dots, s_S\}$$

、ロボットがとりうる行動の集合

$$\mathcal{A} = \{a_1, a_2, \dots, a_A\}$$

、ロボットがある状態である行動を取った際の状態の遷移の仕方を定める遷移関数

$$\mathcal{P} = \{p_{ss'}^a\}_{s,s',a} \text{ (ただし } \sum_{s'} p_{ss'}^a = 1 \text{)}$$

、ロボットがある状態とった行動の良さに関する情報を与える報酬関数

$$\mathcal{R} = \{r_1, r_2, \dots, r_S\}$$

、未来に受け取る報酬の考慮度合いをコントロールする割引率（ただし、 $0 < \gamma < 1$ ）の5つの組

10

20

30

40

$$(S, A, P, R, \gamma)$$

で定義される。

【 0 0 3 8 】

このMDPの設定のもと、ロボットには各状態でどの行動を実行するかの自由度が与えられる。このロボットが各状態 s にいる時に実行する行動 a を決定する関数を方策と呼び、 $\pi(s)$ と書く。ここで、

$$\pi(s) \in A$$

10

であり、 $\pi(s)$ で状態 s にいるときに実行する行動を表す。

【 0 0 3 9 】

強化学習では複数存在する方策のうち、最も現在から将来にいたるまで得られる報酬の期待割引和を最大化する方策、最適方策を求める。最適方策を導く際に重要な役割を果たすのが価値関数 Q である。

【 0 0 4 0 】

【 数 1 】

20

$$Q^\pi(s, a) = \lim_{T \rightarrow \infty} E^\pi \left[\sum_{k=0}^T \gamma^k \mathcal{R}(S_k) \middle| S_0 = s, A_0 = a \right] \quad \dots (1)$$

【 0 0 4 1 】

ここで、 S_k は、ある時刻 k における状態であり、 S_0 は集合

S

における最初の状態 s を表す。

30

【 0 0 4 2 】

この関数は、状態 s で行動 a を実行し、そのあとは方策 π にしたがって無限に行動し続けた場合に得られる報酬の期待割引和を表している。方策 π が最適方策であったとき、最適方策における価値関数 Q^* (最適価値関数) は、

【 数 2 】

$$Q^*(s, a) = \mathcal{R}(s) + \gamma \sum_{s'} p_{ss'}^a \max_{a'} Q^*(s', a') \quad \dots (2)$$

を満たすことが知られ、この式のことをベルマン最適方程式と呼ぶ。

40

Q学習に代表される強化学習の多くの手法は、上記の式の関係性を利用して、この最適価値関数をまずはじめに推定し、その結果を用いて、下記式(3)と設定することで最適方策 π^* を得ている。

【 0 0 4 3 】

【 数 3 】

$$\pi^*(s) = \arg \max_a Q^*(s, a) \quad \dots (3)$$

【 0 0 4 4 】

< < 強化学習による信号制御 > >

50

【 0 0 4 5 】

単独の信号制御を行う既存技術（非特許文献 2）は前節の強化学習のアプローチにもとづき信号の制御策を発見している。信号制御の場合、MDPにおける行動は例えば「南北方向を青、東西方向を赤にする」、「南北方向を赤、東西方向を青にする」という信号の設定を切り替える操作に対応する。

【 0 0 4 6 】

同様に状態は、上記行動の決定の際に利用できる情報、例えば道路に設置してあるセンサや車両から送信される情報などに対応する。

【 0 0 4 7 】

報酬は、例えば「信号で停止している車両の台数の符号反転（-1をかける）」と設定しておけば、停止する車両の数が少なくなるような最適方策が発見されると期待できる。

【 0 0 4 8 】

前述した既存技術（非特許文献 2）は、上記の設定に加え、Space Invader などゲーム分野で大きな成功を収めている Deep Q - Network (DQN) (参考文献 1) のアプローチを採用している。

【 0 0 4 9 】

[参考文献 1]

Human-level control through deep reinforcement learning, Mnih, Volodymyr and Kavukcuoglu, Koray and Silver, David and Rusu, Andrei A and Veness, Joel and Bellemare, Marc G and Graves, Alex and Riedmiller, Martin and Fidjeland, Andreas K and Ostrovski, Georg and others, Nature, 2015.

【 0 0 5 0 】

この研究では、最適価値関数をパラメタをもつニューラルネットワークで近似することを考える。

【 0 0 5 1 】

【 数 4 】

$$Q^*(s, a) \approx Q(s, a; \theta) \quad \cdot \cdot \cdot \quad (4)$$

【 0 0 5 2 】

このパラメタを、シミュレーションを通して学習することで最適価値関数と最適方策を得ている。このアプローチは、とりうる状態の数が非常に多い場合に有効な方法であることが知られている。信号制御の場合でも、センサ数やそこから送られてくる情報の種類が多い場合、一つ一つの状態が数 10 ~ 数 100 次元のベクトルで表現される場合があり、このようなアプローチが採用されている。

【 0 0 5 3 】

このような DQN を用いた強化学習によって、既存技術（参考文献 2）では最適方策によって、車両の待ち時間を減少させることができたと報告されている。

【 0 0 5 4 】

<< 原理 >>

本発明に係る実施形態は、上記既存技術を大幅に発展させ、信号機と車両の移動経路を同時に最適化することによってさらに車両の待ち時間を減少させる技術である。

【 0 0 5 5 】

信号機と車両の移動経路とを同時に最適化を行うために、本実施形態では、移動経路指示機、という仮想的な機械を導入する。

【 0 0 5 6 】

この移動経路指示機は、ある道路の一定区間を走行している車両に対して進むべき経路を指示する、というものである。

【 0 0 5 7 】

例えば、図 4 に示すように、その区間を通行した車両のその区間の通過後にとりうる経

10

20

30

40

50

路としてルート 1 またはルート 2 の 2 種類が存在するとき、この移動経路指示機は、その車両がルート 1 に進むか、ルート 2 に進むべきかを指示する。

【0058】

なお、この指示には、車両が必ず従うとしても良いし、一定の確率に従って指示に従うか否かが決まるとしてもよい。

【0059】

提案エージェントは、信号機と上記移動経路指示機の両方の行動を決定するエージェントである。図 4 の例では、提案エージェントの取りうる行動は全 8 種類であり、各状態においてどの行動を選択するかを定める方策を、強化学習を用いて学習する。

【0060】

上記提案エージェントの行動を反映することのできる交通シミュレータなどを用意することで、提案エージェントの価値関数パラメタ推定には、任意の強化学習手法が適用でき、DQN (参考文献 1) 以外の手法、例えば Double DQN (参考文献 2) などを用いてもよい。

【0061】

[参考文献 2]

Deep Reinforcement Learning with Double Q-Learning, Van Hasselt, Hado and Guez, Arthur and Silver, David, AAAI, 2016.

【0062】

上記の提案エージェントの優れた点は、エージェントのとりうる行動数が車両の数に依存しない、という点にある。

【0063】

車両毎に移動経路を決定する、というエージェントの定義では、アクションの数が、(信号機のアクション数) × (車両 1 のアクション数 (ルート数)) × (車両 2 のアクション数 (ルート数)) ... と車両数の増大にともないアクションの数が指数的に増大していく。

【0064】

すなわち、このような膨大なアクション数を持つエージェントを定義してしまうと、探索空間が増大し、強化学習によって正しく価値関数を推定することが極めて困難となる。

【0065】

提案エージェントでは、このような困難さを回避することで、信号機と車両の移動経路を最適化する最適方策と価値関数を推定する。

【0066】

<本発明の実施の形態に係る交通制御システムの構成>

図 1 を参照して、本発明の実施の形態に係る交通制御システムの構成について説明する。図 1 は、本発明の実施の形態に係る交通制御システムの構成を示すブロック図である。

【0067】

図 1 に示すように、本実施形態に交通制御システム 100 は、価値関数パラメタ学習装置 1 と、外部装置 2 と、センサ 3 と、信号情報指示装置 4 と、移動経路指示装置 5 とを備えて構成される。

【0068】

価値関数パラメタ学習装置 1 は、CPU と、RAM と、後述する行動選択学習処理ルーチンを実行するためのプログラムを記憶した ROM とを備えたコンピュータで構成され、機能的には次に示すように構成されている。

【0069】

価値関数パラメタ学習装置 1 は、シミュレーション設定処理部 10 と、価値関数モデル設定処理部 20 と、価値関数パラメタ推定部 30 と、価値関数制御パラメタ処理部 40 と、記録部 50 と、入出力部 60 とを備えて構成される。

【0070】

シミュレーション設定処理部 10 は、入出力部 60 から取得したシミュレーションを行

10

20

30

40

50

うために必要な情報を、シミュレーション設定記録部 5 1 に格納する。

【 0 0 7 1 】

具体的には、シミュレーション設定処理部 1 0 は、シミュレーション設定に関する情報、例えば、道路ネットワーク、信号位置、センサ位置等をシミュレーション設定記録部 5 1 に格納する。

【 0 0 7 2 】

価値関数モデル設定処理部 2 0 は、入出力部 6 0 から取得した価値関数モデルの設定に関する情報を、価値関数モデル設定記録部 5 2 に格納する。

【 0 0 7 3 】

具体的には、価値関数モデル設定処理部 2 0 は、価値関数モデルの設定に関する情報、例えばニューラルネットワークの層数、中間素子数、活性化関数等を価値関数モデル設定記録部 5 2 に格納する。

10

【 0 0 7 4 】

価値関数パラメタ推定部 3 0 は、シミュレーション設定記録部 5 1 に記録されているシミュレーション設定に関する情報、価値関数モデル設定記録部 5 2 に記録されている価値関数モデルの設定に関する情報を入力とし、同時動的制御シミュレーション実行部 3 1 の処理を繰り返し実行することで、価値関数パラメタを学習し、学習した価値関数パラメタを価値関数パラメタ記録部 5 3 に格納する。

【 0 0 7 5 】

価値関数パラメタ推定部 3 0 は、同時動的制御シミュレーション実行部 3 1 と、価値関数パラメタ学習部 3 2 とを備えて構成される。

20

【 0 0 7 6 】

同時動的制御シミュレーション実行部 3 1 は、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行する。

【 0 0 7 7 】

具体的には、同時動的制御シミュレーション実行部 3 1 は、価値関数パラメタ学習部 3 2 から、シミュレーション設定に関する情報と、価値関数モデルの設定に関する情報とを取得して、シミュレーションを実行する。

30

【 0 0 7 8 】

より具体的には、同時動的制御シミュレーション実行部 3 1 は、シミュレーション設定に関する情報に基づいて、道路に設置してあるセンサから送信される情報を、交通状況を表す状態とし、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を、行動とし、信号で停止している車両の台数の符号反転 (- 1 をかける) を、報酬として、価値関数を用いて、行動を決定し、決定した行動を行ったときの交通状況をシミュレーションする。

【 0 0 7 9 】

ここで計算される報酬は、状態 s から求められる報酬関数

40

$$r_S$$

であるが、状態 s における行動 a も加味した報酬関数

$$r_S(a)$$

としても良い。

【 0 0 8 0 】

そして、同時動的制御シミュレーション実行部 3 1 は、シミュレーションの結果を価値

50

関数パラメタ学習部 3 2 に渡す。

【 0 0 8 1 】

価値関数パラメタ学習部 3 2 は、同時動的制御シミュレーション実行部 3 1 によるシミュレーションの結果に基づいて、価値関数のパラメタを学習する。

【 0 0 8 2 】

具体的には、まず、価値関数パラメタ学習部 3 2 は、シミュレーション設定記録部 5 1 からシミュレーション設定に関する情報を取得し、価値関数モデル設定記録部 5 2 から価値関数モデルの設定に関する情報を取得する。

【 0 0 8 3 】

次に、価値関数パラメタ学習部 3 2 は、シミュレーション設定に関する情報と、価値関数モデルの設定に関する情報とを、同時動的制御シミュレーション実行部 3 1 に渡し、予め定めた反復条件を満たすまで、同時動的制御シミュレーション実行部 3 1 にシミュレーションを繰り返し実行させる。

10

【 0 0 8 4 】

ここで、反復条件は、所定回数を繰り返す、価値関数パラメタに変化が無くなった、価値関数パラメタの学習が収束した等、様々な条件を設定することができる。

【 0 0 8 5 】

そして、価値関数パラメタ学習部 3 2 は、同時動的制御シミュレーション実行部 3 1 から取得したシミュレーション結果に基づいて、価値関数パラメタを学習する。

【 0 0 8 6 】

より具体的には、同時動的制御シミュレーション実行部 3 1 によるシミュレーション結果から得られる報酬と価値関数の値とに基づいて、最適方策を得ることができる価値関数パラメタ（例えば、式（4）におけるパラメタ）を学習する。

20

【 0 0 8 7 】

その後、価値関数パラメタ学習部 3 2 は、学習した価値関数パラメタを、価値関数パラメタ記録部 5 3 に記録する。

【 0 0 8 8 】

価値関数制御パラメタ処理部 4 0 は、価値関数パラメタ記録部 5 3 に記録されている価値関数パラメタを、入出力部 6 0 に渡す。

【 0 0 8 9 】

記録部 5 0 は、シミュレーション設定記録部 5 1 と、価値関数モデル設定記録部 5 2 と、価値関数パラメタ記録部 5 3 とを備えて構成される。

30

【 0 0 9 0 】

シミュレーション設定記録部 5 1 は、シミュレーション設定処理部 1 0 から取得したシミュレーションを行うために必要な情報を記録している。

【 0 0 9 1 】

また、シミュレーション設定記録部 5 1 は、予め設定されたシミュレーションを行うために必要な情報を記録している。

【 0 0 9 2 】

価値関数モデル設定記録部 5 2 は、価値関数モデル設定処理部 2 0 から取得した価値関数モデルの設定に関する情報を記録している。

40

【 0 0 9 3 】

また、価値関数モデル設定記録部 5 2 は、予め設定された価値関数モデルの設定に関する情報を記録している。

【 0 0 9 4 】

価値関数パラメタ記録部 5 3 は、価値関数パラメタ学習部 3 2 により学習された価値関数パラメタを記録している。

【 0 0 9 5 】

入出力部 6 0 は、外部装置 2 から、シミュレーションを行うために必要な情報と、価値関数モデルの設定に関する情報とを受け付ける。

50

【 0 0 9 6 】

入出力部 6 0 は、シミュレーションを行うために必要な情報が入力されると、シミュレーション設定処理部 1 0 に、シミュレーションを行うために必要な情報を渡す。

【 0 0 9 7 】

入出力部 6 0 は、価値関数モデルの設定に関する情報が入力されると、価値関数モデル設定処理部 2 0 に、価値関数モデルの設定に関する情報を渡す。

【 0 0 9 8 】

また、入出力部 6 0 は、価値関数制御パラメタ処理部 4 0 から、価値関数パラメタを受け取ると、外部装置 2 へ出力する。

【 0 0 9 9 】

外部装置 2 は、シミュレーションを行うために必要な情報と、価値関数モデルの設定に関する情報とを設定する装置であり、予め設定されたシミュレーションを行うために必要な情報や予め設定された価値関数モデルの設定に関する情報に修正・変更がある場合に、修正・変更を受け付ける。

10

【 0 1 0 0 】

そして、修正・変更を受け付けると、修正・変更されたシミュレーションを行うために必要な情報および / または価値関数モデルの設定に関する情報を、入出力部 6 0 に渡す。

【 0 1 0 1 】

また、外部装置 2 は、入力された交通状況を表すセンサ情報と、入力された価値関数パラメタとを、信号情報指示装置 4 と、各移動経路指示装置 5 とにそれぞれ渡す。

20

【 0 1 0 2 】

具体的には、まず、外部装置 2 は、入出力部 6 0 から価値関数パラメタと、複数のセンサ 3 の各々から、当該センサ 3 によって計測された交通状況を表すセンサ情報とを取得する。ここで、センサ情報は、車両の速度、車両の台数、車両が通ったか否かなどのそのセンサによって得られる車両の情報である。

【 0 1 0 3 】

次に、外部装置 2 は、信号情報指示装置 4 と、各移動経路指示装置 5 とに対して、取得した交通状況を表すセンサ情報と、価値関数パラメタとを渡す。

【 0 1 0 4 】

センサ 3 は、道路に複数設置されているセンサであり、各設置地点における交通の状況を計測する。例えば、設置地点の画像や設置地点を通過した車両の速度、所定時間内の車両台数、車両が通過したこと等を計測する。

30

【 0 1 0 5 】

信号情報指示装置 4 は、価値関数パラメタ学習装置 1 によって学習された価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する状態について、最適方策となる行動を決定し、決定された行動に含まれる各信号機に対する指示に従って、各信号機に対して指示を行う。

【 0 1 0 6 】

具体的には、まず、信号情報指示装置 4 は、外部装置 2 から、センサ情報と、価値関数パラメタとを取得し、取得したセンサ情報に対応する状態を求める。

40

【 0 1 0 7 】

次に、信号情報指示装置 4 は、取得した価値関数パラメタを用いて、求めた交通状況の状態について、最適方策となる行動を決定する。

【 0 1 0 8 】

そして、信号情報指示装置 4 は、決定された行動に含まれる各信号機に対する指示にしたがって、各信号機に対して、「赤にする」、「青にする」等の指示を行う。

【 0 1 0 9 】

移動経路指示装置 5 は、価値関数パラメタ学習装置 1 によって学習された価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する状態について、最適方策となる行動を決定し、決定された行動に含まれる自装置の区間を通過する移動体に対し

50

て行う進むべき経路の指示に従って、自装置の区間を通過する移動体に対して行う進むべき経路の指示を行う。

【0110】

具体的には、まず、移動経路指示装置5は、外部装置2から、センサ情報と、価値関数パラメタとを取得し、取得したセンサ情報に対応する状態を求める。

【0111】

次に、移動経路指示装置5は、取得した価値関数パラメタを用いて、求めた交通状況の状態について、最適方策となる行動を決定する。

【0112】

そして、移動経路指示装置5は、決定された行動に含まれる自装置の区間を通過する移動体に対して行う進むべき経路の指示にしたがって、自装置の区間を通過する車両に対して「ルート1に進む」、「ルート2に進む」等の指示を行う。

【0113】

<本発明の実施の形態に係る価値関数パラメタ学習装置の作用>

図2は、本発明の実施の形態に係る価値関数パラメタ学習処理ルーチンを示すフローチャートである。

【0114】

価値関数パラメタ学習装置1に価値関数パラメタ学習処理の実行命令がなされると、価値関数パラメタ推定部30において、図2に示す価値関数パラメタ学習処理ルーチンが実行される。

【0115】

まず、ステップS100において、価値関数パラメタ学習部32は、シミュレーション設定記録部51からシミュレーション設定に関する情報を取得する。このシミュレーション設定に関する情報は、予め設定されたものでもよいし、シミュレーション設定処理部10により格納されたものでもよい。

【0116】

ステップS110において、価値関数パラメタ学習部32は、価値関数モデル設定記録部52から価値関数モデルの設定に関する情報を取得する。この価値関数モデルの設定に関する情報は、予め設定されたものでもよいし、価値関数モデル設定処理部20により格納されたものでもよい。

【0117】

ステップS120において、価値関数パラメタ学習部32は、価値関数パラメタを初期化する。

【0118】

ステップS130において、同時動的制御シミュレーション実行部31は、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを実行する。

【0119】

ステップS140において、価値関数パラメタ学習部32は、同時動的制御シミュレーション実行部31によるシミュレーションの結果に基づいて、価値関数のパラメタを学習する。

【0120】

ステップS150において、価値関数パラメタ学習部32は、予め定めた反復条件を満たすか否かを判定する。

【0121】

予め定めた反復条件を満たしていない場合(ステップS150のNO)、ステップS130~S140の処理を繰り返す。

【0122】

10

20

30

40

50

予め定めた反復条件を満たしている場合（ステップS150のYES）、ステップS160において、入出力部60は、ステップS140により学習された価値関数パラメタを、外部装置2へ出力する。

【0123】

<本発明の実施の形態に係る信号情報指示装置4及び移動経路指示装置5の作用>

図3は、本発明の実施の形態に係る制御処理ルーチンを示すフローチャートである。

【0124】

外部装置2から価値関数パラメタが入力されると、信号情報指示装置4において、図3に示す制御処理ルーチンが実行される。

【0125】

まず、ステップS200において、信号情報指示装置4は、外部装置2から入力された、価値関数パラメタ学習装置1によって学習された価値関数パラメタを取得する。

【0126】

ステップS210において、信号情報指示装置4は、外部装置2から、各センサ3のセンサ情報を取得し、取得したセンサ情報に対応する状態を求める。

【0127】

ステップS220において、信号情報指示装置4は、ステップS200で取得した価値関数のパラメタを用いて、ステップS210により求められた交通状況を表すセンサ情報に対応する状態について、最適方策となる行動を決定する。

【0128】

ステップS230において、信号情報指示装置4は、ステップS220により決定した行動に含まれる各信号機に対する指示に従って、各信号機に対して指示を行う。

【0129】

また、各移動経路指示装置5においても、上記図3に示す制御処理ルーチンと同様の処理ルーチンを実行し、決定した行動に含まれる、自装置の区間を通過する移動体に対して行う進むべき経路の指示を含む行動に従って、自装置の区間を通過する移動体に対して行う進むべき経路の指示を行う。

【0130】

以上説明したように、本実施形態に係る価値関数パラメタ学習装置によれば、交通状況を表す状態と、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動との組み合わせに対する価値関数を用いて、各区間を通過する移動体に対して行う進むべき経路の指示、及び各信号機に対する指示を含む行動を行ったときの交通状況のシミュレーションを行い、シミュレーションの結果に基づいて、価値関数パラメタを学習するため、移動体の数が増加しても、最適な交通状況を実現するための価値関数パラメタを学習することができる。

【0131】

また、本実施形態に係る外部装置によれば、価値関数パラメタ学習装置によって学習された価値関数のパラメタを用いて、入力された交通状況を表すセンサ情報に対応する状態について、行動を決定し、決定された行動に従って、各信号機に対する指示、及び各区間を通過する移動体に対して行う進むべき経路の指示を行うため、移動体の数が増加しても、最適な交通状況を実現することができる。

【0132】

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く、この発明の要旨を逸脱しない範囲内で様々な変形や応用が可能である。

【0133】

本実施形態では、信号情報指示装置4と、各移動経路指示装置5とが、各センサ3から取得したセンサ情報に基づいて状態を求め、行動を決定する例を説明したが、外部装置2が、各センサ3からセンサ情報を取得し、それぞれセンサ情報に対応する状態を求めて、行動を決定する構成としても良い。この場合、外部装置2が、信号情報指示装置4と、各移動経路指示装置5に対して、各信号機に対する指示、及び経路の指示を出す。

10

20

30

40

50

【 0 1 3 4 】

また、本願明細書中において、プログラムが予めインストールされている実施形態として説明したが、当該プログラムを、コンピュータ読み取り可能な記録媒体に格納し、またはネットワークを介して提供することも可能である。

【 符号の説明 】

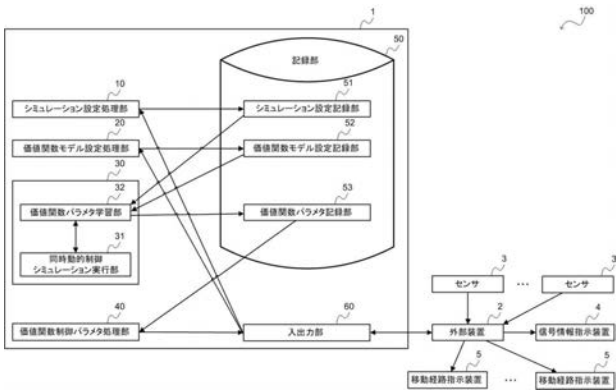
【 0 1 3 5 】

- 1 価値関数パラメタ学習装置
- 2 外部装置
- 3 センサ
- 4 信号情報指示装置
- 5 移動経路指示装置
- 10 シミュレーション設定処理部
- 20 価値関数モデル設定処理部
- 30 価値関数パラメタ推定部
- 31 同時動的制御シミュレーション実行部
- 32 価値関数パラメタ学習部
- 40 価値関数制御パラメタ処理部
- 50 記録部
- 51 シミュレーション設定記録部
- 52 価値関数モデル設定記録部
- 53 価値関数パラメタ記録部
- 60 入出力部
- 100 交通制御システム

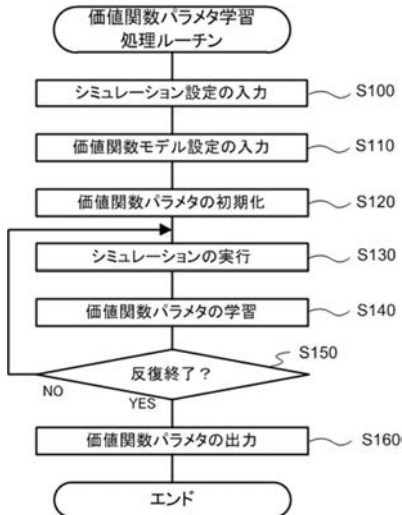
10

20

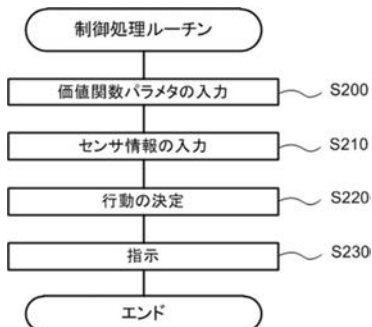
【 図 1 】



【 図 2 】



【 図 3 】



フロントページの続き

(72)発明者 戸田 浩之

東京都千代田区大手町一丁目5番1号 日本電信電話株式会社内

Fターム(参考) 5H181 AA05 AA21 AA26 BB04 DD02 DD03 EE03 FF12 JJ02 JJ06