

(19)日本国特許庁(JP)

(12)特許公報(B2)

(11)特許番号

特許第7036545号

(P7036545)

(45)発行日 令和4年3月15日(2022.3.15)

(24)登録日 令和4年3月7日(2022.3.7)

(51)国際特許分類

F I

G 0 5 B 13/02 (2006.01)

G 0 5 B 13/02

J

請求項の数 19 外国語出願 (全36頁)

(21)出願番号	特願2017-131700(P2017-131700)	(73)特許権者	507342261
(22)出願日	平成29年7月5日(2017.7.5)		トヨタ モーター エンジニアリング ア
(65)公開番号	特開2018-37064(P2018-37064A)		ンド マニファクチャリング ノース
(43)公開日	平成30年3月8日(2018.3.8)		アメリカ, インコーポレイティド
審査請求日	令和2年7月3日(2020.7.3)		アメリカ合衆国、7 5 0 2 4 テキサス
(31)優先権主張番号	15/205,558		州、ブレイノ、ダブリュ 1 - 3 シー・ヘ
(32)優先日	平成28年7月8日(2016.7.8)		ッドクォーターズ・ドライブ、6 5 6 5
(33)優先権主張国・地域又は機関	米国(US)	(74)代理人	100099759
			弁理士 青木 篤
		(74)代理人	100123582
			弁理士 三橋 真二
		(74)代理人	100092624
			弁理士 鶴田 準一
		(74)代理人	100147555
			弁理士 伊藤 公一

最終頁に続く

(54)【発明の名称】 能動的探索なしの強化学習に基づくオンライン学習法及び車両制御方法

(57)【特許請求の範囲】

【請求項 1】

車両の自律的動作を適応的に制御するコンピュータ実行型方法であって、該方法は、

a) 車両を自律的に制御するように構成されたコンピュータ処理システムにおけるcriticネットワークにおいて、受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、actorネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、

b) コンピュータ処理システム内においてcriticネットワークに対して作用的に連結されたactorネットワークにおいて、車両に対して適用されて到達コストに対する最小値を生成する制御入力を決定すること、とを備え、

前記actorネットワークは、推定平均コストと、近似された到達コスト関数から決定された推定到達コストと、車両の現在の状態に対する制御用動力学的值と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成され、

前記近似された到達コスト関数は、以下の関係に従い、重み付けされた放射基底関数の線形結合を用いて決定され、

【数 1】

$$\hat{Z}(x) := \sum_{j=0}^N \omega_j f_j(x)$$

式中、 ω_j は重みであり、 f_j は第 j 番目の放射基底関数であり、 N は近似された到達コスト関数を決定するために使用される放射基底関数の個数であり、且つ、

$\hat{Z}(x)$

10

は近似された到達コスト関数である、方法。

【請求項 2】

前記近似された到達コスト関数において使用される重み ω_j は、以下の関係に従って更新され、

【数 2】

$$\omega^{i+1} = \tilde{\omega}^i + \lambda_1 + \sum_{l=0}^N \lambda_{2l} \delta_{jl} + \lambda_3 \hat{Z}_{avg} f_k$$

20

式中、 δ_{ij} はディラックのデルタ関数を表し、上付き文字 i は反復の回数を表し、 λ_1 、 λ_2 、 λ_3 はラグランジュ乗数であり、且つ

\hat{Z}_{avg}

は推定平均コストである、請求項 1 に記載の方法。

30

【請求項 3】

ベルマン方程式の線形化版を用いて決定された近似的な時間差的誤差を用いて前記 critic ネットワークのパラメータを更新する段階を更に備える、請求項 1 に記載の方法。

【請求項 4】

前記 critic ネットワークのパラメータを更新する段階は、前記車両が運動しているときに実施される、請求項 3 に記載の方法。

【請求項 5】

車両の自律的動作を適応的に制御するコンピュータ実行型方法であって、該方法は、

a) 車両を自律的に制御するように構成されたコンピュータ処理システムにおける critic ネットワークにおいて、受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、actor ネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、

40

b) コンピュータ処理システム内において critic ネットワークに対して作用的に連結された actor ネットワークにおいて、車両に対して適用されて到達コストに対する最小値を生成する制御入力を決定すること、とを備え、

前記 actor ネットワークは、推定平均コストと、近似された到達コスト関数から決定された推定到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成され、

当該方法は、ベルマン方程式の線形化版を用いて決定された近似的な時間差的誤差を用い

50

て前記criticネットワークのパラメータを更新する段階を更に備え、

前記criticネットワークにより決定された推定平均コストは、以下の関係に従って更新され、

【数 3】

$$\hat{Z}_{avg}^{i+1} = \hat{Z}_{avg}^i - \alpha_2^i e_k \hat{Z}_k$$

10

式中、 α_2^i は学習率であり、 e_k は近似的な時間差的誤差であり、

$$\hat{Z}_k$$

は前記近似された到達コスト関数から決定された推定コストであり、

$$\hat{Z}_{avg}^i$$

は状態 i における推定平均コストであり、且つ、

20

$$\hat{Z}_{avg}^{i+1}$$

は状態 $i + 1$ における推定平均コストである、方法。

【請求項 6】

前記受動的に収集されたデータは、前記criticネットワークのパラメータを更新する間に使用される唯一のデータである、請求項 3 に記載の方法。

【請求項 7】

車両の自律的動作を適応的に制御するコンピュータ実行型方法であって、該方法は、

a) 車両を自律的に制御するように構成されたコンピュータ処理システムにおけるcriticネットワークにおいて、受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、actorネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、

30

b) コンピュータ処理システム内においてcriticネットワークに対して作用的に連結されたactorネットワークにおいて、車両に対して適用されて到達コストに対する最小値を生成する制御入力を決定すること、とを備え、

前記actorネットワークは、推定平均コストと、近似された到達コスト関数から決定された推定到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成され、

当該方法は、以下の関係に従って決定された近似的な時間差的誤差を用いて前記criticネットワークのパラメータを更新する段階を更に備え、

40

【数 4】

$$e_k := \hat{Z}_{avg} \hat{Z}_k - \exp(-q_k) \hat{Z}_{k+1}$$

式中、 e_k は近似的な時間差的誤差であり、

50

\hat{Z}_{avg}

は推定平均コストであり、

\hat{Z}_k

は状態 k における推定到達コストであり、

\hat{Z}_{k+1}

10

は状態 $k + 1$ における推定到達コストであり、且つ、 q_k は状態 k における状態コストである、方法。

【請求項 8】

前記近似された到達コスト関数は、前記 critic ネットワークにおいてリアルタイムで学習される、請求項 1 に記載の方法。

【請求項 9】

車両の自律的動作を適応的に制御するコンピュータ実行型方法であって、該方法は、

a) 車両を自律的に制御するように構成されたコンピュータ処理システムにおける critic ネットワークにおいて、受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、actor ネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、

20

b) コンピュータ処理システム内において critic ネットワークに対して作用的に連結された actor ネットワークにおいて、車両に対して適用されて到達コストに対する最小値を生成する制御入力を決定すること、とを備え、

前記 actor ネットワークは、推定平均コストと、近似された到達コスト関数から決定された推定到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成され、

前記ノイズレベルは、以下の関係に従い、重み付けされた基底関数の線形結合を用いて学習され、

30

【数 5】

$$\rho(x; \mu) \approx \hat{\rho}(x; \mu) := \sum_j^M \mu_j g_j(x)$$

式中、 μ は推定ノイズレベルであり、 μ_j は、 g_j により表された第 j 番目の放射基底関数に対する重みであり、且つ、 M は、ノイズレベルを推定するために使用されるべき放射基底関数の個数である、方法。

40

【請求項 10】

以下の関係に従って決定された近似誤差を用いて前記 actor ネットワークの重み付けパラメータを更新する段階を更に備え、

【数 6】

$$d_k \approx q_k \Delta t - \hat{V}_{k+1} + \hat{V}_k + \hat{V}_{avg} + L_{k,k+1} \rho_k,$$

式中、 d_k は近似誤差であり、 q_k は状態 k における状態コストであり、

\hat{V}_k

10

は状態 k において近似された到達コストであり、

\hat{V}_{k+1}

は状態 $k + 1$ において近似された到達コストであり、

\hat{V}_{avg}

は近似された平均コストであり、且つ、

【数 7】

20

$$L_{k,k+1} := (0.5\hat{V}_k - \hat{V}_{k+1})^\top B_k B_k^\top \hat{V}_k \Delta t$$

であり、式中、 B_k は状態 k における制御用動力学的値である、請求項 9 に記載の方法。

【請求項 1 1】

前記actorネットワークの重み付けパラメータの更新は、前記車両が運動しているときに実施される、請求項 1 0 に記載の方法。

30

【請求項 1 2】

前記actorネットワークの重み付けパラメータは、以下の関係に従って更新され、

【数 8】

$$\mu^{i+1} = \mu^i - \beta^i d_k L_{k,k+1} g_k,$$

式中、

μ^{i+1}

40

は状態 $i + 1$ における重み付けパラメータの値であり、

μ^i

は状態 i における重み付けパラメータの値であり、 β^i は学習率であり、 d_k は時間差的誤差であり、且つ g は放射基底関数である、請求項 1 0 に記載の方法。

【請求項 1 3】

50

前記制御入力を用い、前記自律的動作を制御すべく使用可能な制御ポリシーを修正する段階を更に備える、請求項 1 に記載の方法。

【請求項 14】

前記推定平均コストが収束するまで、前記段階 (a) 及び (b) を反復的に実施して前記制御入力を再決定することにより、前記自律的動作を制御するために使用可能な制御ポリシーを最適化する段階を更に備える、請求項 1 に記載の方法。

【請求項 15】

前記制御ポリシーは、能動的探索なしで最適化される、請求項 14 に記載の方法。

【請求項 16】

車両の自律的動作を適応的に制御するように構成されたコンピュータ処理システムであって、該コンピュータ処理システムは、該コンピュータ処理システムの動作を制御する一つ以上のプロセッサと、該一つ以上のプロセッサにより使用可能なデータ及びプログラム命令を記憶するメモリとを備え、

前記一つ以上のプロセッサは、前記メモリ内に記憶された命令を実行して、

a) 受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、前記車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定し、且つ、

b) 前記車両に対して適用されて前記到達コストに対する最小値を生成する制御入力を決定する、ように構成され、

前記一つ以上のプロセッサは、前記推定平均コストと、前記近似された到達コスト関数から決定された到達コストと、前記車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成され、

前記近似された到達コスト関数は、以下の関係に従い、重み付けされた放射基底関数の線形結合を用いて決定され、

【数 1】

$$\hat{Z}(x) := \sum_{j=0}^N \omega_j f_j(x)$$

式中、 ω_j は重みであり、 f_j は第 j 番目の放射基底関数であり、N は近似された到達コスト関数を決定するために使用される放射基底関数の個数であり、且つ、

$\hat{Z}(x)$

は近似された到達コスト関数である、コンピュータ処理システム。

【請求項 17】

前記一つ以上のプロセッサは、前記メモリ内に記憶された命令を実行し、前記推定平均コストが収束するまで、前記段階 (a) 及び (b) を反復的に実施して前記制御入力を再決定することにより、前記自律的動作を制御するために使用可能な制御ポリシーを最適化するように構成される、請求項 16 に記載のコンピュータ処理システム。

【請求項 18】

コンピュータシステムにより実行可能な命令が自身内に記憶された、一時的でないコンピュータ可読媒体であって、

a) 受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、

b) 前記車両に対して適用されて前記到達コストに対する最小値を生成する制御入力を決
定すること、とを備える機能を実施させ、

前記制御入力は、前記到達コストに対する最小値を生成し、且つ前記推定平均コストと、
前記近似された到達コスト関数から決定された到達コストと、前記車両の現在の状態に対
する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベ
ルを推定することにより制御入力決定され、

前記近似された到達コスト関数は、以下の関係に従い、重み付けされた放射基底関数の線
形結合を用いて決定され、

【数 1】

$$\hat{Z}(x) := \sum_{j=0}^N \omega_j f_j(x)$$

式中、 ω_j は重みであり、 f_j は第 j 番目の放射基底関数であり、N は近似された到達コスト
関数を決定するために使用される放射基底関数の個数であり、且つ、

$\hat{Z}(x)$

は近似された到達コスト関数である、一時的でないコンピュータ可読媒体。

【請求項 19】

前記命令は、前記推定平均コストが収束するまで、前記段階 (a) 及び (b) を反復的に
繰り返して前記制御入力を再決定することにより、自律的動作を制御するために使用可能
な制御ポリシーを最適化するように実行可能である、請求項 18 に記載の一時的でないコ
ンピュータ可読媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、車両を自律的に制御する方法に関し、更に詳細には、車両の操作を自律的に制
御するために使用可能な制御ポリシー (control policy) を修正及び / 又は最適化するた
めの強化学習方法に関する。

【背景技術】

【0002】

一定形式のシステムにおいては、周囲環境を能動的に探索することにより、最適なシステ
ム制御ポリシーを決定するために、モデルフリー (model-free) 強化学習 (RL) 技術が
採用され得る。しかし、車両が採用し得る全ての動作の膨大な能動的探索 (active explo
ration) に伴う潜在的に否定的な結果により、車両の自律的制御に対して使用可能な制御
ポリシーに対して従来の RL 手法を適用することは困難であり得る。これに加え、車両安
全性の確保を支援するために必要とされる態様で能動的探索を行うと、大きなコンピュー
タ処理コストが必要とされ得る。代替策として、車両が動作している周囲環境の正確な動
力学的モデルを利用することにより、能動的探索なしで最適な制御ポリシーを決定すべく
、モデルベースの RL 技術が採用され得る。しかし、自律車両が動作している複雑な周囲
環境は、正確にモデル化することが非常に困難なことがある。

【発明の概要】

【0003】

本明細書中に記述された実施形態の一つの見地においては、車両の自律的動作を適応的に
(adaptively) 制御するコンピュータ実行型方法が提供される。該方法は、(a) 車両を
自律的に制御すべく構成されたコンピュータ処理システムにおける critic ネットワークに

において、受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、actorネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定する段階と、(b) コンピュータ処理システム内においてcriticネットワークに対して作用的に連結されたactorネットワークにおいて、車両に対して適用されて到達コストに対する最小値を生成する制御入力を決定する段階とを含み、actorネットワークは、平均コストと、近似された到達コスト関数から決定された到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータとを用いて、ノイズレベルを推定することにより制御入力を決定すべく構成される。

【0004】

本明細書中に記述された実施形態の別の見地においては、車両の自律的動作を適応的に制御するように構成されたコンピュータ処理システムが提供される。該コンピュータ処理システムは、該コンピュータ処理システムの動作を制御する一つ以上のプロセッサと、該一つ以上のプロセッサにより使用可能なデータ及びプログラム命令を記憶するメモリとを含み、上記一つ以上のプロセッサは、メモリ内に記憶された命令を実行して、(a) 受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定し、且つ(b) 車両に対して適用されて到達コストに対する最小値を生成する制御入力を決定する、ように構成され、一つ以上のプロセッサは、平均コストと、到達コスト関数から決定された到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成される。

【0005】

本明細書中に記述された実施形態の別の見地においては、一時的でないコンピュータ可読媒体が提供される。該媒体は、コンピュータシステムにより実行可能な命令を該媒体内に記憶し、該コンピュータシステムに、受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、車両に対して適用されて到達コストに対する最小値を生成する制御入力を決定することとを備える機能を実施させ、制御入力は、到達コストに対する最小値を生成し、且つ、平均コストと、到達コスト関数から決定された到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力決定される。

【図面の簡単な説明】

【0006】

【図1】本明細書中に記述された実施形態に係る、(例えば自律車両などの)システムに対する制御入力を決定すべく且つシステム制御ポリシーを修正及び/又は最適化すべく構成されたコンピュータ処理システムのブロック図である。

【図2】本明細書中に記述された方法に係る、車両制御入力の決定、及び/又は、制御ポリシーの修正若しくは最適化の間における情報の流れを示す概略図である。

【図3】制御入力を決定し且つ制御ポリシーを修正及び/又は最適化する方法の実施形態の動作を示すフローチャートである。

【図4】本明細書中に記述された実施形態に係る、一つ以上の制御入力と制御ポリシーとを使用する自律的制御に向けて構成された車両であって、当該車両に対する制御入力を決定すべく且つ自律車両操作制御ポリシーを修正及び/又は最適化すべく構成されたコンピュータ処理システムが組み込まれた車両の概略的ブロック図である。

【図5】本明細書中に記述された実施形態に係る方法を用いる、高速道路合流用の制御ポリシーの最適化の例において採用された車両の構成の概略図である。

【図6】図5に示された車両の構成に関して実施される最適化のグラフ表示である。

【発明を実施するための形態】

【0007】

本明細書中に記述された実施形態は、コンピュータ実行型の強化学習(RL)方法に関す

10

20

30

40

50

るものであり、この強化学習方法は、車両を自律的に制御すべく使用可能な制御入力を決
定するため、及び車両の操作を自律的に制御する制御ポリシーを修正及び／又は最適化す
べく制御入力を使用するために使用可能である。この方法は、（例えば、動作を実施し、
且つその動作の結果を監視して制御ポリシーを決定及び改変することを伴い得る）能動的
探索を使用せずに、制御入力を決定し且つ制御ポリシーを最適化し得る。本明細書中に記
述された方法は、能動的探索の代わりに、受動的に収集されたデータと、部分的に既知で
あるシステムの動力学的モデルと、制御されている車両に関する既知の制御用動力学的モ
デルとを使用する。

【 0 0 0 8 】

本開示に関連して、「オンライン」とは、コンピュータ処理システムが学習し得ると共に
、actor及びcriticのネットワークパラメータが、上記システムが動作するにつれて（例え
ば車両が移動するなどにつれて）、コンピュータ処理され且つ更新され得ることを意味す
る。オンラインのソリューションを用いてactorパラメータ及びcriticパラメータを決定かつ
更新すると、車両及びシステムの動力学的値（dynamics）の変更が許容され得る。同様
に、自律的動作とは、自律的に実施される動作である。

【 0 0 0 9 】

図 1 は、本明細書中に開示される種々の実施形態に係る方法を実現すべく構成されたコン
ピュータ処理システム 1 4 のブロック図である。更に詳細には、少なくとも一つの実施形
態において、コンピュータ処理システム 1 4 は、本明細書中に記述された方法に従い、制
御入力を決定すべく構成され得る。また、コンピュータ処理システムは、システム（例え
ば、自律車両）を制御して特定の操作若しくは機能を自律的に実施すべく使用可能な制御
ポリシーを修正及び／又は最適化するようにも構成され得る。

【 0 0 1 0 】

少なくとも一つの実施形態において、コンピュータ処理システムは、車両に組み込まれ得
ると共に、生成された制御入力を使用して車両の操作の制御に向けられた制御ポリシーを
修正及び／又は最適化すべく構成され得る。制御入力を決定するため及び制御ポリシーを
修正及び／又は最適化するためにコンピュータ処理システムにより必要とされる情報（例
えば、データ、命令、及び／又は他の情報）の少なくとも幾つかは、任意の適切な手段か
ら、例えば車両センサから又は無線接続を介して遠隔データベースのような車外情報源か
ら、受信され且つ／又はそれにより収集され得る。幾つかの実施形態においては、制御ポ
リシーを修正及び／又は最適化するためにコンピュータ処理システムにより必要とされる
情報（例えば、データ）の少なくとも幾つかは、車両の操作の前に（例えば、メモリ内に
記憶されたデータ及び他の情報として）コンピュータ処理システムに提供され得る。また
、コンピュータ処理システムは、制御入力に従って且つ／又は修正若しくは最適化された
制御ポリシーに従って車両を制御することで、関連する自律的動作を実施するようにも構
成され得る。

【 0 0 1 1 】

少なくとも一つの実施形態において、コンピュータ処理システムは、（例えばスタンドア
ロンのコンピュータ処理システムとして）車両から遠隔的に配置され得ると共に、制御入
力を決定すべく且つ車両の自律的動作の実施に向けられた制御ポリシーを修正及び／又は
最適化するように構成され得る。遠隔的なコンピュータ処理システムによって生成された
最適化又は修正された制御ポリシーは、車両による展開のために車両のコンピュータ処理
システムへロード又はインストールされて、実際の交通環境において車両を制御し得る。

【 0 0 1 2 】

図 1 を参照すると、コンピュータ処理システム 1 4 は、コンピュータ処理システム 1 4 及
び関連する構成要素の全体的な操作を制御する（少なくとも一つのマイクロプロセッサを
含み得る）一つ以上のプロセッサ 5 8 であって、メモリ 5 4 のような一時的でない（non-
transitory）コンピュータ可読媒体内に記憶された命令を実行する、プロセッサ 5 8 を含
み得る。本開示に関連して、コンピュータ可読記憶媒体とは、命令を実行するシステム、
装置若しくはデバイスによって使用されるか又はそれに関連して使用されるプログラムを

10

20

30

40

50

含む又は記憶し得る任意の有形媒体であり得る。プロセッサ 58 は、プログラムコード中に含まれた命令を実施すべく構成された少なくとも一つのハードウェア回路（例えば、集積回路）を含み得る。複数のプロセッサ 58 が在る構成において、斯かるプロセッサは相互から独立して作動し得るか、又は、一つ以上のプロセッサが相互に協働して作動し得る。

【0013】

幾つかの実施形態において、コンピュータ処理システム 14 は、RAM 50、ROM 52、及び/又は他の任意で適切な形態のコンピュータ可読メモリを含み得る。メモリ 54 は、一つ以上のコンピュータ可読メモリを備え得る。一つ又は複数のメモリ 54 は、コンピュータ処理システム 14 の構成要素であり得るか、又は、一つ又は複数のメモリは、コンピュータ処理システム 14 に作用的に接続されてコンピュータ処理システム 14 に使用され得る。本説明を通して使用される「作用的に接続された」という語句は、直接的な物理接触のない接続を含め、直接的又は間接的な接続を含み得る。

10

【0014】

一つ以上の構成において、本明細書中に記述されたコンピュータ処理システム 14 は、人工的又はコンピュータ的な知能要素、例えば、ニューラルネットワーク、ファジィ論理回路、又は他の機械学習アルゴリズム、を組み込み得る。更に、一つ以上の構成において、本明細書中に記述された特定の機能又は操作を実施するように構成されたハードウェア及び/又はソフトウェア要素は、複数の要素及び/又は箇所に分散され得る。コンピュータ処理システム 14 に加え、車両は、コンピュータ処理システム 14 により実施される制御機能を増強若しくは支援するために、又は他の目的のために、付加的なコンピュータ処理システム及び/又はデバイス（図示せず）を組み込み得る。

20

【0015】

メモリ 54 は、データ 60、及び/又は、プロセッサ 58 によって実行されて種々の機能を実行し得る命令（例えば、プログラムロジック）56 を含み得る。データ 60 は、受動的に収集されたデータを含み得る。受動的に収集されたデータは、能動的探索から収集されたのではないデータとして定義され得る。受動的に収集されたデータの一例は、ビルの頂部に取付けられたカメラを用いて高速道路の入口の周りにおける車両の軌跡の獲得を記述する、<http://www.fhwa.dot.gov/publications/research/operations/06137/> に記述されたデータセットである。別の例において、受動的に収集されたデータは、人間の運転者により実行された操縦に応じて車両センサにより収集されたデータを含み得る。人間の運転者により実行された操縦と、この操縦が実行された車両環境条件と、この操縦に引き続き且つ/又はこの操縦に応じて車両周囲において生じた事象と、に関して、データが収集されてコンピュータ処理システムに提供され得る。或いは、コンピュータ処理システムが車両に設置されたときに、コンピュータ処理システム 14 は、（制御ポリシー 101 のような）一つ以上の車両制御ポリシーのオンラインでの修正及び/又は最適化のために、斯かる受動的に収集されたデータを収集及び/又は受信するように構成され得る。

30

【0016】

車両制御用動力学的モデル 87 は、種々の入力に対して車両が如何に応答するかを記述する刺激応答モデル（stimulus-response model）であり得る。車両制御用動力学的モデル 87 は、本明細書中に記述されるように、車両に対する制御入力を決定する上で且つ制御ポリシー 101 を修正及び/又は最適化する上で使用されるべく、（受動的に収集されたデータを用いて）所定の車両状態 x における車両に対して状態コスト $q(x)$ 及び制御用動力学的値 $B(x)$ を決定するように使用され得る。与えられた任意の車両に対する車両制御用動力学的モデル 87 が決定されて、メモリ内に記憶され得る。

40

【0017】

再び図 1 を参照すると、コンピュータ処理システムの実施形態は、2つの学習システム又は学習ネットワーク、並びに相互に作用するactorネットワーク（又は「actor」）83 及びcriticネットワーク（又は「critic」）81 も含み得る。これらネットワークは、例えば、人工ニューラルネットワーク（ANN）を用いて実現され得る。

【0018】

50

本明細書中に記述された目的に対し、(変数 によっても表される)制御ポリシー 101 は、一群の車両の状態のうちの各状態 x に応じて車両により取られるべき動作 u を特定又は決定する関数又は他の関係として定義され得る。故に、自律的動作の実行中の車両の各状態 x に対し、車両は、関連する動作 $u = (x)$ を実施するように制御され得る。したがって、制御ポリシーは、車両の操作を制御して、例えば、高速道路合流などの関連する操作を自律的に実施する。actor 83 は、制御ポリシーに関して動作し、critic から受信した情報及び他の情報を用いて、ポリシーを修正及び/又は最適化し得る。制御ポリシーにより自律的に制御された車両操作は、高速道路への合流、又は、車線の変更のような特定の目的を達成すべく実施される一つの運転操作又は一群の運転操作として定義され得る。

【0019】

コンピュータ処理システム 14 は、制御ポリシーの修正及び最適化に対して使用可能である新規な半モデルフリー RL 方法 (semi-model-free RL method) (本明細書においては受動的 actor/critic (pAC) 方法という) を実行するように構成される。この方法において、critic は、車両の種々の状態に対する評価関数を学習し、且つ、actor は、能動的探索なしで、代わりに受動的に収集されたデータと既知の車両制御用動力学的モデルとを用いて制御ポリシーを改善する。この方法は、部分的に既知であるシステムの動力学的モデルを使用することにより、能動的探索に対する必要性を回避する。この方法は、車両環境の制御されていない動力学的値又は過渡的なノイズレベルに関する知見を必要としない。この方法は、例えば、環境がノイズ的に如何に展開するかのサンプルは入手可能であるが車両センサにより能動的に探索することは困難であり得る自律車両に関して、実行可能である。

【0020】

制御入力を決定し且つ制御ポリシーを修正及び/又は最適化する目的に対し、状態 $x \in \mathbb{R}^n$ 及び制御入力 $u \in \mathbb{R}^m$ により、離散時間確率論的動力学系は以下のように定義され得る。

【数 1】

$$\Delta x = A(x_t)\Delta t + B(x_t)u_t\Delta t + C(x_t)d\omega \quad (1)$$

式中、 (t) はブラウニアン運動であり、

$A(x_t)$

、

$B(x_t)u_t$

及び

$C(x_t)$

は、それぞれ、受動的動力学的値、車両制御用動力学的値、及び、過渡的ノイズレベルである。 Δt は、時間のステップサイズである。この種の系は、多くの状況において生ずる (例えば、ほとんどの機械系のモデルはこれらの動力学に従う)。関数 $A(x)$ 、 $B(x)$ 及び $C(x)$ は、理解されるべく、モデル化されている特定の系に依存する。受動的動力学的値は、車両の環境における変化であって、車両システムに対する制御入力の結果ではない変化を含む。

【0021】

本明細書中に記述された方法及びシステムにおいて、離散時間動力学系に対する MDP は、タプル

$$\langle \mathcal{X}, \mathcal{U}, P, R \rangle$$

であり、式中、

$$\mathcal{X} \subseteq \mathbb{R}^n$$

及び

$$\mathcal{U} \subseteq \mathbb{R}^m$$

10

は、状態空間及び動作空間である。

$$P := \{p(\mathbf{y}|\mathbf{x}, \mathbf{u}) | \mathbf{x}, \mathbf{y} \in \mathcal{X}, \mathbf{u} \in \mathcal{U}\}$$

は、動作による状態遷移モデルであり、且つ、

$$R := \{r(\mathbf{x}, \mathbf{u}) | \mathbf{x} \in \mathcal{X}, \mathbf{u} \in \mathcal{U}\}$$

20

は、状態 \mathbf{x} 及び動作 \mathbf{u} に関する即時コスト関数である。先に記述されたように、制御ポリシー

$$\mathbf{u} = \pi(\mathbf{x})$$

は、状態 \mathbf{x} から動作

$$\mathbf{u}$$

30

へとマッピングする関数である。予期される累積コストである、ポリシー の下での到達コスト関数 (cost-to-go function) (又は価値関数)

$$V^\pi(\mathbf{x})$$

は、無限時間区間 (infinite horizon) の平均コストの最適性判断基準の下で、以下のよう

【数 2】

$$V^\pi(x_t) := \sum_{k=1}^{\infty} p(\mathbf{x}_k, \pi(\mathbf{x}_k)) r(\mathbf{x}_k, \pi(\mathbf{x}_k)) \Delta t - V_{avg}^\pi$$

40

$$V_{avg}^\pi := \lim_{N \rightarrow \infty} \frac{1}{N \Delta t} \sum_{k=1}^N p(\mathbf{x}_k, \pi(\mathbf{x}_k)) r(\mathbf{x}_k, \pi(\mathbf{x}_k)) \Delta t$$

50

式中、

$$V_{avg}^{\pi}$$

は平均コストであり、 k は時間インデックスであり、且つ、 t は時間ステップである。
最適な到達コスト関数は、以下の離散時間ハミルトン - ヤコビ - ベルマン方程式を満足する。

【数 3】

$$V_{avg}^{\pi} + V^{\pi}(\mathbf{x}_k) = \min_{\mathbf{u}_k} Q^{\pi}(\mathbf{x}_k, \mathbf{u}_k) \quad (2)$$

式中、 $Q^{\pi}(\mathbf{x}_k, \mathbf{u}_k) := r(\mathbf{x}_k, \mathbf{u}_k)\Delta t + G[V^{\pi}](\mathbf{x}_k)$, であり、且つ、
 $G[V^{\pi}](\mathbf{x}_k) := \int_{\mathbf{X}} p(\mathbf{y}|\mathbf{x}_k, \pi) V^{\pi}(\mathbf{y}) d\mathbf{y}$ である。

式中、

$$Q^{\pi}(\mathbf{x}, \mathbf{u})$$

は動作価値関数であり、且つ、

$$G[\cdot]$$

は積分演算子である。MDPの目的は、以下の関係に従い、無限時間区間に亘り、平均コストを最小化する制御ポリシーを見出すことである。

【数 4】

$$\pi^*(\mathbf{x}_k) = \arg \min_{\pi} \mathbb{E}[V_{avg}^{\pi}]$$

ここで、最適な制御ポリシーにおける値は、上付き文字 $*$ を以て表され得る（例えば、

$$V^*$$

、

$$V_{avg}^*$$

）。

【0022】

離散時間動力学系に対する線形MDP（L-MDP）は、連続的な状態空間及び動作空間に対して厳密な解が迅速に求められ得るという利点を備えた汎用マルコフ決定過程のサブクラスである。構築された動学的値、及び、別体的な状態コスト及び制御コストの下で、ベルマン方程式は、組み合わされた状態コスト及び制御されていない動学的値の線形

10

20

30

40

50

固有関数を見出すことに解が限定された線形微分方程式として再構築され得る。その後、L - M D P に対する到達コスト関数（又は、価値関数）は、正確な動力学モデルが利用可能であるときに、二次プログラミング（Q P）のような最適化方法により、効率的に求められ得る。

【 0 0 2 3 】

マルコフ決定過程の線形公式は、以下に示されるように、制御コストを定義すべく、且つ、車両動力学値に関する条件を加えるべく使用され得る。

【数 5】

$$r(\mathbf{x}_k, \mathbf{u}_k) := q(\mathbf{x}_k) + KL(p(\mathbf{x}_{k+1}|\mathbf{x}_k) \| p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k)) \quad (3)$$

$$p(\mathbf{x}_{k+1}|\mathbf{x}_k) = 0 \Rightarrow \forall \mathbf{u}_k, p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k) = 0 \quad (4)$$

ここで、

$$q(\mathbf{x}) \geq 0$$

は状態コスト関数であり、且つ、

$$KL(\cdot \| \cdot)$$

はクルバック - ライブラー（K L）偏差である。式（3）は、動作のコストを、それが系に対して有する確率論的效果の量に対して関連付け、且つ、それを状態コストに対して加算する。第2の条件は、何らの動作も、受動的動力学の下では達成され得ない新たな遷移を導入しないことを確実にする。式（1）により表された確率論的動力学系は、当然、上記仮定を満足する。

【 0 0 2 4 】

ハミルトン - ヤコビ - ベルマン方程式（式（2））は、L - M D P 形態において、指数的に変換された到達コスト関数に対する線形微分方程式（以下、線形化ベルマン方程式という）へと書き換えられ得る。

【数 6】

$$Z_{avg} Z(\mathbf{x}_k) = \exp(-q(\mathbf{x}_k) \Delta t) \mathcal{G}[Z](\mathbf{x}_{k+1})$$

$$Z(\mathbf{x}) := \exp(-V^*(\mathbf{x})), \quad (5)$$

$$Z_{avg} := \exp(-V_{avg}^*)$$

$$p(\mathbf{x}_{k+1}|\mathbf{x}_k, \pi_k^*) = \frac{p(\mathbf{x}_{k+1}|\mathbf{x}_k) Z(\mathbf{x}_{k+1})}{\mathcal{G}[Z](\mathbf{x}_{k+1})}$$

式中、

$$Z(\mathbf{x})$$

及び Z_{avg} は、それぞれ、 Z 値と称される指数的に変換された到達コスト関数、及び、最適ポリシーの下での平均コストである。(式(1))における状態遷移はガウス性であることから、制御された動力学的値と受動的な動力学的値との間のKL偏差は、

【数7】

$$KL(p(\mathbf{x}_{k+1}|\mathbf{x}_k) \parallel p(\mathbf{x}_{k+1}|\mathbf{x}_k, \mathbf{u}_k)) = \frac{1}{2\rho(\mathbf{x}_k)} \mathbf{u}_k^\top \mathbf{u}_k$$

$$\frac{1}{\rho(\mathbf{x}_k)} := B(\mathbf{x}_k)^\top (C(\mathbf{x}_k)^\top C(\mathbf{x}_k))^{-1} B(\mathbf{x}_k) \quad (6)$$

として表される。

【0025】

その後、L-MDP系に対する最適な制御ポリシーは、

【数8】

$$\pi^*(\mathbf{x}_k) = -\rho(\mathbf{x}_k) B(\mathbf{x}_k)^\top V_{\mathbf{x}_k} \quad (7)$$

として表され、式中、

$V_{\mathbf{x}_k}$

は、 \mathbf{x}_k における \mathbf{x} に関する到達コスト関数 V の偏微分値である。 Z 値及び平均コストは、系の動力学的値が完全に入手可能であるとき、固有値又は固有関数を解くことにより、線形化ベルマン方程式から導かれ得る。

【0026】

本明細書中に記述されたコンピュータ処理システム14の実施形態は、種々の形式の入力及び出力情報を測定、受信、及び/又は、アクセスすることによりシステム(例えば、車両)の状態 $\mathbf{x}(t)$ を決定する。例えば、データは、このシステムに結合されたセンサ又はこのシステムと別途通信するセンサを用いて測定され得る。コンピュータ処理システム14は、制御入力 \mathbf{u} を決定することで、式(1)により特徴付けられる車両の安定性及び所望の運動を達成し且つ式(2)において記述されたエネルギーに基づくコスト関数を最小化する。

【0027】

本明細書中に記述されたコンピュータ処理システム14の実施形態は、相互に作用する2つの学習システム又は学習ネットワーク、すなわちactorネットワーク(又は「actor」)83及びcriticネットワーク(又は「critic」)81を含む。これらネットワークは、人工ニューラルネットワーク(ANN)を用いて実現され得る。actor83は、状態依存の制御ポリシーを使用して、車両に対して適用され且つ到達コスト(cost-to-go)に対する最小値を生成する制御入力 $\mathbf{u}(\mathbf{x})$ を決定する。actorは、平均コストと、近似的到達コスト関数から決定された推定到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータとを用いてノイズレベルを推定することにより、制御入力を決定すべく構成される。critic81は、受動的に収集されたデータのサンプル及び状態コ

ストを用いて、推定平均コストと、actorネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定する。本明細書中に開示される幾つかの実施形態において、actor 8 3 は内部ループフィードバックコントローラとして実現されると共に、critic 8 1 は、外部ループフィードバックコントローラとして実現される。両者ともに、制御命令をもたらすように作動可能である車両を起動可能な機構又は制御器に関するフィードフォワード経路中に配置される。

【 0 0 2 8 】

受動的に収集されたデータのサンプルと、車両制御用動力学的モデル 8 7 から受信された状態コスト $q(x)$ とを用いて、critic 8 1 は、車両の現在の状態 x_k 、次の状態 x_{k+1} 、及び、最適ポリシー下での状態コスト q_k を評価し、且つ、先に記述されたベルマン方程式の線形化版（式（ 5 ））を使用して、近似された到達コスト関数

$$\hat{Z}(x)$$

（ Z 値 ）を決定し、且つ、actor 8 3 により使用される推定平均コスト

$$\hat{Z}_{avg}$$

を生成する。Z 値を推定するために、重み付けされた放射基底関数（ R B F ）の線形結合（ linear combination ）が使用され得る。

【 数 9 】

$$Z(x) \approx \hat{Z}(x) := \sum_{j=0}^N \omega_j f_j(x)$$

式中、 ω_j は重みであり、 f_j は第 j 番目の R B F であり、且つ、N は R B F の個数である。基底関数は、車両システムの非線形の動力学に依存して適切に選択され得る。重みは、指数化された真の到達コストと推定到達コストとの間における最小二乗誤差を最小化することにより、最適化される。

$$Z(x_k)$$

及び

$$Z_{avg}$$

を真の Z 値及び真の平均 Z 値コストとし、且つ

$$\hat{Z}_k$$

、

$$\hat{Z}_{avg}$$

をそれぞれそれらの推定対応物とする。近似された到達コスト関数は、critic ネットワークによりリアルタイムで学習され得る。

【数 1 0】

$$\min_{\omega, \hat{Z}_{avg}} \frac{1}{2} \sum_D (\hat{Z}_{avg} \hat{Z}_k - Z_{avg} Z_k)^2, \quad (8)$$

$$s.t. \sum_{i=0}^N \omega_i = C, \forall i \omega_i \geq 0, \forall \mathbf{x} \hat{Z}_{avg} \hat{Z}(\mathbf{x}) \leq 1, \quad (9)$$

10

式中、Cは自明な解 $= 0$ への収束を回避すべく使用される一定値である。第2及び第3の制約は、式(5)から由来する

$$\forall \mathbf{x}, 0 < Z_{avg} Z(\mathbf{x}) \leq 1,$$

と、

$$\forall \mathbf{x}, q(\mathbf{x}) \geq 0.$$

20

とを満足するために必要とされる。

【0 0 2 9】

重み 及び平均コスト

 \hat{Z}_{avg}

は、真の到達コストと推定到達コストとの間の誤差を、線形化ベルマン方程式(LBE)(式(5))から以下のように決定された近似的な時間差的誤差 e_k により近似することにより、ラグランジュ緩和時間差(TD)学習に基づいて更新され得る。なぜなら、真の到達コスト及び真の平均コストは、pAC方法に対して使用された情報によっては決定されないからである。

30

【数 1 1】

$$\begin{aligned} \hat{Z}_{avg} \hat{Z} - Z_{avg} Z_k &\approx e_k := \hat{Z}_{avg} \hat{Z}_k - \exp(-q_k) \hat{Z}_{k+1} \\ \tilde{\omega}^{i+1} &= \omega^i - \alpha_1^i e_k Z_{avg} \mathbf{f}_k \end{aligned} \quad (10)$$

$$\omega^{i+1} = \tilde{\omega}^i + \lambda_1 + \sum_{l=0}^N \lambda_{2l} \delta_{jl} + \lambda_3 \hat{Z}_{avg} \mathbf{f}_k \quad (11)$$

40

$$\hat{Z}_{avg}^{i+1} = \hat{Z}_{avg}^i - \alpha_2^i e_k \hat{Z}_k \quad (11A)$$

式中、

 α_1^i

及び

50

α_2

は、学習率であり、且つ、 e_k はL-MDPに対するTD誤差である。 δ_{ij} はディラックのデルタ関数を表している。下付き文字*i*は、反復の回数を表している。 λ_1 、 λ_2 、 λ_3 は、制約式(9)に対するラグランジュ乗数である。 λ_1 は、式(10)による誤差を最小化すべく、且つ、式(11)による制約を満足すべく更新される。反復とは、(criticに対する重み w 及びactorに対する μ のような)critic及びactorのパラメータの更新として定義され得る。これに加え、criticネットワークのパラメータの更新は、車両が運動しているときに実施され得る。本明細書中に記述された方法において、criticネットワークの更新の間に使用される唯一のデータは、受動的に収集されたデータである。

10

【0030】

各乗数の値は、以下の方程式を解くことにより算出される。

【数12】

$$\begin{bmatrix} \sum_j f_j & f_0 \cdots f_N & f^T f \\ 1 & 1 \cdots 0 & f_0 \\ \vdots & \ddots & \vdots \\ 1 & 0 \cdots 1 & f_N \\ N & 1 \cdots 1 & \sum_j f_j \end{bmatrix} \begin{bmatrix} \lambda_1^i \\ \lambda_2^i \\ \lambda_3^i \end{bmatrix} = \begin{bmatrix} c - \sum_j \tilde{w}^i \\ -\tilde{w}_0^i \\ \vdots \\ -\tilde{w}_N^i \\ 1 - \hat{Z}_{avg} \hat{Z}_k \end{bmatrix}$$

20

【0031】

幾つかの場合、制約の部分集合は有効でないことがあり得る。斯かる場合、これらの制約に対する乗数はゼロに設定され、且つ、残りの有効な制約に対する乗数が求められる。criticは、受動的動力学の下での状態遷移サンプル

$(\mathbf{x}_k, \mathbf{x}_{k+1})$

30

及び状態コスト q_k を用いて、各パラメータを更新する。重み w 、推定Z値

\hat{Z}

、及び、平均コスト

\hat{Z}_{avg}

40

は、車両が運動している間に、式(10)～(11A)に従い、オンラインで更新され得る。

【0032】

コンピュータ処理システムにおいて、criticネットワークに作用的に結合されたactor 83は、制御入力を決定し、到達コストに対する最小値を生成する車両に適用され得る。criticにより生成された推定到達コスト

\hat{Z}

及び推定平均コスト

50

\hat{Z}_{avg}

と、状態コスト $q(x)$ と、車両制御用動力学的モデル 87 から決定された現在の状態に対する制御用動力学的値情報 $B(x)$ と、criticにより使用されて到達コスト関数

 $\hat{Z}(x)$

を推定且つ推定平均コスト

 \hat{Z}_{avg}

10

を生成すべく受動的に収集されたデータのサンプルとを用い、actor 83 は制御入力決定し得る。制御入力は、制御ポリシーを修正すべく使用され得る。特定の実施形態において、ポリシーは、上述の態様で、収束するまで反復的に修正され、その時点でそれは最適化されたと見做される。actorは、標準的なベルマン方程式を用い、且つ、能動的探索なしで、制御ポリシーを改善する。制御用動力学的値は、車両に対する既知の制御用動力学的モデルから決定され得る。

【0033】

actor 83 はまた、制御入力 $u(x)$ を各車両システムに対してリアルタイムで適用し、所望の操作（例えば、高速道路合流、車線変更）を自律的に実施もし得る。本明細書中に開示される幾つかの実施形態において、actor 83 は、内部ループフィードバックコントローラにおいて具現され得ると共に、critic 81 は、外部ループフィードバックコントローラにおいて具現され得る。両者ともに、車両を起動し得る制御器に関するフィードフォワード経路中に配置される。

20

【0034】

actor 83 は、criticからの評価値（例えば、

 \hat{Z}

30

及び

 \hat{Z}_{avg}

）、受動的動力学の下でのサンプル、及び、既知の制御用動力学的値を用いて、ノイズレベルを推定することにより、制御ポリシーを改善又は修正する。ノイズレベルは、重み付けされた各放射基底関数の線形結合により、近似的に学習される。

【数13】

40

$$\rho(x; \mu) \approx \hat{\rho}(x; \mu) := \sum_j^M \mu_j g_j(x) \quad (12)$$

式中、 μ_j は、 j 番目の放射基底関数 g_j に対する重みである。 M は、放射基底関数の個数である。

【0035】

は、到達コストと動作 - 状態値との間の最小二乗誤差を最小化することにより、最適化される。

50

【数 1 4】

$$\min_{\mu} \frac{1}{2} \sum_D (\hat{Q}_k - V_k^* - V_{avg}^*)^2$$

式中、

V^*

10

、

V_{avg}^*

及び

\hat{Q}

20

は、最適な制御ポリシーの下での、真の到達コスト関数、平均コスト、及び推定された動作 - 状態値である。最適制御ポリシーは目的関数を最小化することにより学習され得る。なぜなら、真の動作 - 価値コストは、ポリシーが最適ポリシーであるとき且つそのときにのみ、

$V^* + V_{avg}^*$

に等しいからである。

\hat{Z}_k

30

及び

Z_{avg}^*

は、ノイズレベルを更新するときに、以下の関係に従い、

\hat{V}_k

及び

40

\hat{V}_{avg}

を決定すべく使用され得る。

【数 1 5】

$$\hat{V}_k = -\log(\hat{Z}_k)$$

$$\hat{V}_{avg} = -\log(\hat{Z}_{avg})$$

【 0 0 3 6 】

重み μ は、以下に定義される近似的な時間差的誤差 d_k により更新される。標準的なベルマン方程式は近似されて誤差 d_k を決定する。なぜなら、真の到達コスト及び真の平均コストは算出されることができないからである。

10

【 数 1 6 】

$$\hat{Q}_k - V_k^* - V_{avg}^* \approx d_k := \hat{Q}_k - \hat{V}_k - \hat{V}_{avg}$$

$$\hat{Q}_k \approx q_k \Delta t + \frac{0.5 \Delta t}{\hat{\rho}_k} \mathbf{u}_k^\top \mathbf{u}_k + \hat{V}(\mathbf{x}_{k+1} + B_k \mathbf{u}_k \Delta t)$$

$$\mathbf{u}_k \approx -\hat{\rho}_k B_k^\top \hat{V}_k,$$

20

式中、

$$\mathbf{x}_{k+1}$$

は、受動的動力学の下での次の状態であり、且つ、

$$\mathbf{x}_{k+1} + B_k \mathbf{u}_k \Delta t$$

30

は、動作 \mathbf{u}_k による制御された動力学の下での次の状態である。推定到達コスト、平均コスト、及び、それらの微分値は、criticからの推定されたZ値及び平均Z値コストを利用することにより算出され得る。更に、

$$\hat{V}(\mathbf{x}_{k+1} + B_k \mathbf{u}_k \Delta t)$$

は、

【 数 1 7 】

40

$$\hat{V}(\mathbf{x}_{k+1} + B_k \mathbf{u}_k \Delta t) \approx \hat{V}_{k+1} + \hat{V}_{x_{k+1}}^\top B_k \mathbf{u}_k \Delta t.$$

により近似されて、

μ

に関してTD誤差を線形化し得る。

【 0 0 3 7 】

50

μ

は、近似されたTD誤差による時間差(TD)学習を用いて更新され得る。

【数18】

$$\mu^{i+1} = \mu^i - \beta^i d_k L_{k,k+1} g_k,$$

$$d_k \approx q_k \Delta t - \hat{V}_{k+1} + \hat{V}_k + \hat{V}_{avg} + L_{k,k+1} \rho_k,$$

$$L_{k,k+1} := (0.5 \hat{V}_k - \hat{V}_{k+1})^\top B_k B_k^\top \hat{V}_k \Delta t$$

10

式中、 β^i は学習率であり、且つ、 $L_{k,k+1}$ は、項 $L(x_k, x_{k+1})$ の省略版である。

【0038】

この手順は、与えられた状態において、受動的動力学値、状態コスト q_k 、及び、制御用動力学値 B_k の下で、状態遷移サンプル

(x_k, x_{k+1})

20

を用いることにより、能動的探索なしでポリシーを改善する。標準的なactor-critic方法は、能動的探索によりポリシーを最適化する。定義されたこれらのactor及びcriticの機能により、コンピュータ処理システム14は、L-MDPを用いて、半モデルフリー強化学習を実現し得る。

【0039】

本明細書中に記述された方法において、ポリシーは、受動的に収集されたデータのサンプルと、車両制御用動力学の知見とを用い、到達コストと動作-状態値との間の誤差を最小化することにより学習されるパラメータにより最適化される。本明細書中に記述された方法は、乗用車を制御すべく通常的に利用可能である車両自体の動力学モデルにより、最適ポリシーが決定されることを可能とする。上記方法はまた、それらの動力学モデルが通常は既知でない周囲の車両の操作に関して受動的に収集されたデータも使用する。これに加え、本明細書中に記述された方法を用いると、最適な制御ポリシーを決定する上で、車両環境の受動的動力学値 $A(x_t)$ 及び過渡的ノイズレベル $C(x_t)$ は、認識される必要はない。

30

【0040】

図2は、本明細書中に記述された方法に係る、コンピュータ処理システム14における、制御入力の決定、及び、制御ポリシーの修正又は最適化の実行中の情報の流れを示す概略図である。従来のactor-critic方法は、周囲環境から能動的に収集されたデータのサンプルを用いて動作し得る一方、本明細書中に記述されたpAC方法は、周囲環境の能動的探索なしで、代わりに、受動的に収集されたサンプル、及び、既知の車両制御用動力学モデルを用いて、最適な制御ポリシーを決定する。critic81又はactor83において受信された一切の情報は、後で使用するためにメモリ内にバッファリングされ得る。例えば、パラメータ値を算出し又は推定すべくcritic又はactorに必要とされる情報の全てが現在入手できないという状況において、受信情報は、残りの必要な情報が受信されるまで、バッファリングされ得る。

40

【0041】

図3は、本明細書中に開示された幾つかの実施形態に従い、制御入力を決定し且つ制御ポリシーを修正及び/又は最適化するための図1のコンピュータ処理システムの動作を示すフローチャートである。

50

【 0 0 4 2 】

プロセスは、ブロック 3 1 0 にて開始され、そこで critic 8 1 は、推定平均コスト \hat{Z}_{avg}

と、actor ネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数

$$\hat{Z}(x)$$

10

とを決定し得る。

【 0 0 4 3 】

次に、ブロック 3 2 0 にて、actor 8 3 は、到達コスト関数

$$\hat{Z}(x)$$

を用いて、車両に適用されて該車両の到達コストに対する最小値を生成する制御入力決定し得る。actor 8 3 は、制御ポリシーを修正して、このポリシーを改善し且つ / 又はこの制御ポリシーを最適化し得る。

20

【 0 0 4 4 】

ブロック 3 3 0 においては、ブロック 3 2 0 において導かれた制御入力が車両に適用されて、例えば高速道路への合流又は車線の変更などの、車両の自律的動作が行われ得る。また、車両は、任意の改善又は最適化された制御ポリシーに従って更に制御され得る。特定の実施形態において、車両操作は、制御ポリシーが、未だ最適化されたと考えられる点まで改善されていないとしても、ポリシーの最新版を用いてコンピュータ処理システムにより制御され得る。

【 0 0 4 5 】

ブロック 3 4 0 においては、actor 8 3 及び critic 8 1 の種々のパラメータが更新され得る。この更新は、本明細書中に記述された関係に従って実施され得る。この更新に対して使用される唯一のデータは、受動的に収集されたデータであり得る。或る実施形態において、actor 及び critic は、それらのそれぞれのパラメータの更新を実施し得る。或いは、actor 及び critic のパラメータの更新は、ポリシー反復器（図示せず）又は同様の手段により実施され得る。

30

【 0 0 4 6 】

図 4 は、図 1 のコンピュータ処理システム 1 4 と同様の態様で構成されたコンピュータ処理システム 1 1 4 が組み込まれた例示的な実施形態に係る車両 1 1 を示す機能的ブロック図である。車両 1 1 は、乗用車、トラック、又は、本明細書中に記述された操作を実施し得る他の任意の車両の形態を取り得る。車両 1 1 は、完全に又は部分的に自律モードで動作すべく構成され得る。自律モードで動作している間、車両 1 1 は、人的相互作用なしで動作すべく構成され得る。例えば、高速道路の合流操作が実行されている自律モードにおいて、車両は、車両乗員からの入力なしで、高速道路上の車両から安全距離を維持すること、他の車両と速度を調和すること等を行うように、スロットル、ブレーキ及び他のシステムを動作させ得る。

40

【 0 0 4 7 】

車両 1 1 は、コンピュータ処理システム 1 1 4 に加え、且つ、相互に作用的に通信する種々のシステム、サブシステム及び構成要素、及び構成要素、例えば、センサシステム又は配列 2 8、一つ以上の通信インタフェース 1 6、操舵システム 1 8、スロットルシステム 2 0、制動システム 2 2、電源 3 0、動力システム 2 6、並びに本明細書中に記述されたように車両を動作させるために必要な他のシステム及び構成要素を含み得る。車両 1 1 は、図 4 に示されたよりも多い又は少ないサブシステムを含み得ると共に、各サブシステム

50

は、複数の要素を含み得る。更に、車両 11 のサブシステム及び要素の各々は、相互接続され得る。車両 11 の記述された機能及び / 又は自律的動作の一つ以上の実施は、相互に協働して動作している複数の車両システム及び / 又は構成要素により実行され得る。

【 0 0 4 8 】

センサシステム 28 は、任意の適切な形式のセンサを含み得る。本明細書中には、異なる形式のセンサの種々の例が記述される。しかし、実施形態は、記述された特定のセンサに限定されないことは理解される。

【 0 0 4 9 】

センサシステム 28 は、車両 11 の外部環境に関する情報を検知すべく構成された所定数のセンサを含み得る。例えば、センサシステム 28 は、全地球測位システム (GPS) のようなナビゲーションユニット、及び、例えば、慣性測定装置 (IMU) (図示せず)、RADAR ユニット (図示せず)、レーザ測距計 / LIDAR ユニット (図示せず)、及び車両の内部及び / 又は該車両 11 の外部環境の複数の画像を捕捉すべく構成されたデバイスを備える一台以上のカメラ (図示せず) 等の他のセンサを含み得る。カメラは、スチルカメラ又はビデオカメラであり得る。IMU は、慣性加速度に基づいて車両 11 の位置及び向きの変化を検知するように構成されたセンサ (例えば、加速度計及びジャイロスコープ等) の任意の組合せを組み込み得る。例えば、IMU は、車両のロール速度、ヨーレート、ピッチ速度、長手方向加速度、横方向加速度、及び、垂直加速度のようなパラメータを検知し得る。ナビゲーションユニットは、車両 11 の地理的位置を推定すべく構成された任意のセンサであり得る。この目的の為に、ナビゲーションユニットは、地球に対する車両 11 の位置に関する情報を提供するように作動可能な送受信機を含む一つ以上の送受信機を含み得る。また、ナビゲーションユニットは、業界公知の態様で、記憶され且つ / 又は利用可能な地図を用いて与えられた開始点 (例えば、車両の現在位置) から、選択された目的地までの走行ルートを決

10

20

【 0 0 5 0 】

公知の態様において、車両センサ 28 は、種々の車両システムに対する適切な制御命令を策定且つ実行する際にコンピュータ処理システム 114 により使用されるデータを提供する。例えば、慣性センサ、車輪速度センサ、道路状態センサ、及び操舵角センサからのデータは、車両を旋回させるための命令を策定して操舵システム 18 において実行する上で、処理され得る。各車両センサ 28 は、車両 11 に組み込まれる任意の運転者支援機能及び自律的動作機能をサポートするために必要とされる任意のセンサを含み得る。センサシステム 28 が複数のセンサを含む構成において、センサは、相互から独立的に作動し得る。代替的に、各センサのうちの 2 つ以上が、相互に協働して作動し得る。センサシステム 28 のセンサは、コンピュータ処理システム 114 に対し、及び / 又は車両 11 の他の任意の要素に対し、作用的に接続され得る。

30

【 0 0 5 1 】

また、各車両センサ 28 により収集された任意のデータは、本明細書中に記述された目的でデータを必要とし又は利用する任意の車両システム又は構成要素にも送信され得る。例えば、車両センサ 28 により収集されたデータは、コンピュータ処理システム 114 に、又は一つ以上の専用のシステム又は構成要素のコントローラ (図示せず) に送信され得る。付加的な特定の形式のセンサとしては、本明細書中に記述された機能及び操作を実施するために必要とされる他の任意の形式のセンサが挙げられる。

40

【 0 0 5 2 】

特定の車両センサからの情報は、一つよりも多い車両システム又は構成要素を制御すべく処理かつ使用され得る。例えば、自動化された操舵制御及び制動制御の両方を組み込んだ車両において、種々の道路状態センサは、データをコンピュータ処理システム 114 に提供し、このコンピュータ処理システムは、プロセッサが実行可能な記憶された命令に従って道路状態情報を処理すると共に、操舵システム及び制動システムの両方に対して適切な制御命令を策定することができるようになる。

【 0 0 5 3 】

50

車両 11 は、センサの出力信号又は他の信号が、コンピュータ処理システム 114 又は別の車両システム若しくは要素による使用の前に前処理を必要とするという状況、又はコンピュータ処理システムから送信された制御信号が、起動可能なサブシステム又はサブシステム構成要素（例えば、操舵システム又はスロットルシステムの構成要素）による使用の前に処理を必要とするという状況に適した、信号処理手段 38 を含み得る。信号処理手段は、例えば、アナログ／デジタル（A／D）変換器又はデジタル／アナログ（D／A）変換器であり得る。

【0054】

センサ統合機能（sensor fusion capability）138 は、センサシステム 28 からのデータを入力として受け入れるべく構成されたアルゴリズム（又は、アルゴリズムを記憶するコンピュータプログラム製品）の形態であり得る。上記データは、例えば、センサシステム 28 の各センサにて検知された情報を表すデータを含む。センサ統合アルゴリズムは、センサシステムから受信したデータを処理し、（例えば、複数の個別的なセンサの出力から形成された）統合された又は合成された信号を生成し得る。センサ統合アルゴリズム 138 は、例えば、カルマンフィルタ、ベイジアンネットワーク、又は、別のアルゴリズムを含む。センサ統合アルゴリズム 138 は更に、センサシステム 28 からのデータに基づく種々のアセスメントを提供し得る。例示的な実施形態において、アセスメントは、車両 11 の環境における個別的な物体又は特定構造の評価、特定状況の評価、及び、特定の状況に基づく可能的な影響の評価を含み得る。他のアセスメントも可能である。センサ統合アルゴリズム 138 は、コンピュータ処理システム 114 に組み込まれた又はコンピュータ処理システム 114 と作用的に通信する（メモリ 54 のような）メモリ内に記憶され得ると共に、当業界において公知の態様でコンピュータ処理システムにより実行され得る。

【0055】

本明細書中に記述された任意の情報若しくはパラメータの受信、収集、監視、処理、及び／又は、決定を参照するときにおける「連続的に」という語句の使用は、コンピュータ処理システム 114 が、これらのパラメータに関する情報が存在し又は検出されるや否や、又は、センサの取得サイクル及びプロセッサの処理サイクルに従ってできるだけ素早く、任意の情報を受信及び／又は処理すべく構成されることを意味している。コンピュータ処理システム 114 が、例えば、センサからのデータ又は車両構成要素の状況に関する情報を受信すると直ちに、コンピュータ処理システムは、記憶されたプログラム命令に従って動作し得る。同様に、コンピュータ処理システム 114 は、センサシステム 28 から及び他の情報源から、同時進行的又は連続的に情報の流れを受信して処理し得る。この情報は、本明細書中に記述された態様及び目的にて、メモリ内に記憶された命令に従って処理及び／又は評価される。

【0056】

また、図 4 は、先に記述されたように、図 1 のコンピュータ処理システム 14 と同様の態様で構成された代表的なコンピュータ処理システム 114 のブロック図も示している。本明細書中に記述されたようにポリシーの修正を実施すると共に制御入力を決定するために必要とされる機能を組み込むと共に、コンピュータ処理システム 114 は、他の車両システム及び要素に作用的に接続されると共に、その他の点では、車両 11 及びその構成要素の制御及び動作に影響するように構成され得る。コンピュータ処理システム 114 は、少なくとも幾つかのシステム及び／又は構成要素を、（ユーザ入力なしで）自律的に且つ／又は（一定程度のユーザ入力を以て）半自律的に制御すべく構成され得る。また、コンピュータ処理システムは、幾つかの機能を自律的及び／又は半自律的に制御及び／又は実行するようにも構成され得る。コンピュータ処理システム 114 は、種々のサブシステム（例えば、動力システム 26、センサシステム 28、操舵システム 18）から、各通信インタフェース 16 のうちの任意のものから、及び／又は他の任意で適切な情報源から受信した入力及び／又は情報に基づき、車両 11 の機能性を制御し得る。

【0057】

図 4 の実施形態において、コンピュータ処理システム 114 は、図 1 に関して先に記述さ

10

20

30

40

50

れたように、車両制御用動力学的モデル 1 8 7、critic 1 8 1、actor 1 8 3、及び、制御ポリシー 2 0 1 を含み得る。コンピュータ処理システム 1 1 4 は、先に記述されたように、制御入力を決定すべく、且つ自律車両の操作制御ポリシーを修正及び／又は最適化すべく構成され得る。また、コンピュータ処理システム 1 1 4 は、制御入力に従って、且つ、本明細書中に記述されたように修正又は最適化された制御ポリシーにも従って、車両を制御して所望操作を実施すべく構成され得る。

【 0 0 5 8 】

コンピュータ処理システム 1 1 4 は、図 4 に示された要素の幾つか又は全てを有し得る。加えて、コンピュータ処理システム 1 1 4 は、特定の用途に必要とされ又は所望される付加的な構成要素も含み得る。また、コンピュータ処理システム 1 1 4 は、複数のコントローラ又はコンピュータ処理デバイスであって、分散態様にて、情報を処理し且つ／又は車両 1 1 の個別的な構成要素若しくはサブシステムを制御するように機能する複数のコントローラ又はコンピュータ処理デバイスを表し、又は、それにより具現され得る。

10

【 0 0 5 9 】

メモリ 5 4 は、単一又は複数のプロセッサ 5 8 により実行されて、図 1 に関して上述されたものを含む、車両 1 1 の種々の機能を実行するデータ 6 0 及び／又は命令 5 6（例えば、プログラムロジック）を収納し得る。メモリ 5 4 は、本明細書中に記述された車両システム及び／又は構成要素（例えば、動力システム 2 6、センサシステム 2 8、コンピュータ処理システム 1 1 4、及び、通信インタフェース 1 6）のうちの一つ以上にデータを送信し、それらからデータを受信し、それらと相互作用し、又はそれらを制御するための命令を含む、付加的な命令も含み得る。命令 5 6 に加え、メモリ 5 4 は、他の情報の中でも、道路地図、経路情報のようなデータを記憶し得る。斯かる情報は、自律的、半自律的、及び／又は手動的なモードにおける車両 1 1 の動作の間において、ルートを計画するのに且つその他にことをするのに、車両 1 1 及びコンピュータ処理システム 1 1 4 により使用され得る。

20

【 0 0 6 0 】

コンピュータ処理システム 1 1 4 は、（概略的に 6 2 と表される）一つ以上の自律的な機能又は動作を実施するために、種々の起動可能な車両システム及び構成要素の制御を連携調整するように構成され得る。これらの自律的な機能 6 2 は、メモリ 5 4 及び／又は他のメモリ内に記憶されると共に、プロセッサにより実行されたときに、本明細書中に記述された種々のプロセス、命令又は機能のうちの一つ以上を実現するコンピュータ可読プログラムコードの形態で実現され得る。

30

【 0 0 6 1 】

通信インタフェース 1 6 は、車両 1 1 と、外部センサ、他の車両、他のコンピュータシステム、（本明細書中に記述されたように、衛星システム、携帯電話／無線通信システム、種々の車両サービスセンターなどのような）種々の外部のメッセージ及び通信システム、及び／又はユーザとの間の相互作用を許容すべく構成され得る。通信インタフェース 1 6 は、車両 1 1 のユーザに情報を提供し又はユーザから入力を受信するためのユーザインタフェース（例えば、一台以上のディスプレイ（図示せず）、音声／オーディオインタフェース（図示せず）、及び／又は他のインタフェース）を含み得る。

40

【 0 0 6 2 】

また、通信インタフェース 1 6 は、ワイドエリアネットワーク（WAN）、無線通信ネットワーク、及び／又は他の任意で適切な通信ネットワークにおける通信を可能とするインタフェースも含み得る。通信ネットワークは、有線の通信リンク、及び／又は無線の通信リンクを含み得る。通信ネットワークは、上記のネットワーク及び／又は他の形式のネットワークの任意の組合せを含み得る。通信ネットワークは、一つ以上のルータ、スイッチ、アクセスポイント、無線アクセスポイント、及び／又は類似物を含み得る。一つ以上の構成において、通信ネットワークは、任意の近傍車両及び車両 1 1 と、任意の近傍の路側の通信ノード及び／又はインフラとの間の通信を許容し得る、車両対全て（V2X）（車両対インフラストラクチャ（V2I）技術及び車両対車両（V2V）技術を含む）の技術

50

を包含し得る。

【 0 0 6 3 】

W A Nネットワーク環境において使用されたとき、コンピュータ処理システム 1 1 4 は、ネットワーク（例えば、インターネット）のようなW A N上での通信を確立するためのモデム又は他の手段を含み（又は、それに対して作用的に接続され）得る。無線通信ネットワークにおいて使用されたとき、コンピュータ処理システム 1 1 4 は、無線ネットワークにおける一つ以上のネットワークデバイス（例えば、基地送受信ステーション）を介して無線コンピュータ処理デバイス（図示せず）と通信するための一つ以上の送受信機、デジタル信号プロセッサ、及び付加的な回路機構並びにソフトウェアを含み（又は、それに対して作用的に接続され）得る。これらの構成は、種々の外部情報源から定常的な情報の流れを受信する種々の態様を提供する。

10

【 0 0 6 4 】

車両 1 1 は、コンピュータ処理システム 1 1 4 並びに他の車両システム及び / 又は構成要素と作用的に通信し且つコンピュータ処理システムから受信した制御命令に応じて作用し得る、種々の起動可能なサブシステム及び要素を含み得る。種々の起動可能なサブシステム及び要素は、（例えば、A C C 及び / 又は車線維持などの）いずれの自律的の走行支援システムが起動されているのか且つ / 又は車両が完全自律モードで駆動されているのかといった所定の走行状況のような要因に依存して、手動的又は（コンピュータ処理システム 1 1 4 により）自動的に制御され得る。

【 0 0 6 5 】

20

操舵システム 1 8 は、車両ホイール、ラック及びピニオン操舵ギア、操舵ナックル、及び / 若しくは車両 1 1 の方向を調節すべく作用可能であり得る他の任意の要素（コンピュータシステムで制御可能な任意の機構又は要素を含む）、又は要素の組み合わせを含み得る。動力システム 2 6 は、車両 1 1 に動力運動を提供すべく作用可能な構成要素を含み得る。例示的な実施形態において、動力システム 2 6 は、エンジン（図示せず）、（ガソリン、ディーゼル燃料、又は、ハイブリッド車両の場合には一つ以上の電気バッテリーのような）エネルギー源、及び、変速機（図示せず）を含み得る。制動システム 2 2 は、車両 1 1 を減速すべく構成された、要素及び / 又はコンピュータシステムで制御可能な任意の機構の任意の組合せを含み得る。スロットルシステムは、（例えば、加速ペダル、及び / 又は例えばエンジンの作動速度を制御することで車両 1 1 の速度を制御するように構成された任意のコンピュータシステム制御可能な機構などの）要素及び / 又は機構を含み得る。図 1 は、車両に組み込まれ得る車両サブシステムの僅かな例 1 8、2 0、2 2、2 6 を示している。特定の車両は、これらのシステムの一つ以上、又は示されたシステムの一つ以上に加えて他のシステム（図示せず）の一つ以上を組み込み得る。

30

【 0 0 6 6 】

車両 1 1 は、コンピュータ処理システム 1 1 4、センサシステム 2 8、起動可能なサブシステム 1 8、2 0、2 2、2 6、及び他のシステム並びにこれらの要素が、コントローラエリアネットワーク（C A N）バス 3 3 又はその類似物を用いて相互に通信し得るように構成され得る。C A Nバス及び / 又は他の有線又は無線の機構を介し、コンピュータ処理システム 1 4 は、種々の車両システム及び構成要素に対してメッセージを送信し（且つ / 又は、それらからメッセージを受信し）得る。或いは、本明細書中に記述された要素及び / 又はシステムの任意のものは、バスを使用せずに相互に対して直接的に接続され得る。同様に、本明細書中に記述された要素及び / 又はシステム間の接続は、（有線接続のような）別の物理的媒体を経由され得るか、又は上記接続は無線接続であり得る。

40

【 0 0 6 7 】

図 1 は、コンピュータ処理システム 1 4、メモリ 5 4、及び通信インタフェース 1 6 のような車両 1 1 の種々の構成要素を、車両 1 1 に一体化されているとして示しているが、これら構成要素の一つ以上は、車両 1 1 とは別体的に取付けられ、又は関連付けられ得る。例えば、メモリ 5 4 は、部分的に又は完全に、車両 1 1 とは別体的に存在し得る。したがって、車両 1 1 は、別体的又は一体的に配置され得る複数のデバイス要素の形態で提供さ

50

れ得る。車両 1 1 を構成するデバイス要素は、有線又は無線の態様で相互に通信的に結合され得る。

【 0 0 6 8 】

実施例

図 5 及び図 6 を参照すると、本明細書中に記述された制御入力及びポリシー修正 / 最適化方法の実施形態の一つの実施例において、自律的な高速道路合流操作がシミュレートされる。この操作は、4 次元の状態空間、及び、1 次元の動作空間を有する。動力学的値は

【 数 1 9 】

$$\begin{aligned} \mathbf{x} &= [dx_{12}, dv_{12}, dx_{02}, dv_{02}]^T, \\ A(\mathbf{x}) &= [dx_{12}, 0, dv_{02}, +0.5\alpha_0(\mathbf{x})\Delta t, \alpha_0(\mathbf{x})]^T \\ B(\mathbf{x}) &= [0.5\Delta t, 1, 0, 0]^T, C(\mathbf{x}) = [0, 2.5, 0, 2.5]^T \\ \alpha_0(\mathbf{x}) &= \alpha \frac{v_2^\beta (-dv_{02})}{-dx_{02}^\gamma}, \Delta t = 0.1[\text{sec}] \end{aligned}$$

であり、式中、下付き文字 0 は、高速道路の最右側車線上で合流車両の後方の車両（「後続車両」という）を表し、1 は、ランプ R R 上で合流している自動化車両を表し、且つ、2 は、高速道路上の最右側車線上で合流車両の前方の車両（「先行車両」という）を表している。d x₁₂ 及び d v₁₂ は、先行車両からの合流車両の相対的な位置及び速度を表している。例示な目的で、先行車両は一定速度 v₂ = 30 メートル / 秒で走行されること、及び、後続車両に対する車両制御用動力学的モデルは既知であることが仮定される。もし後続車両の速度が先行車両よりも低速である（d v₀₂ < 0）場合には、 $\alpha = 1.55$ 、 $\beta = 1.08$ 、 $\gamma = 1.65$ であり、その他の場合には、 $\alpha = 2.15$ 、 $\beta = -1.65$ 、 $\gamma = -0.89$ である。状態コスト

$q(\mathbf{x})$

は、

【 数 2 0 】

$$q(\mathbf{x}) = k_1 \left(1.0 - \exp \left(-k_2 \left(1 - \frac{2dx_{12}}{dx_{02}} \right)^2 - k_3 (dv_{10})^2 \right) \right)$$

であり、式中、k₁、k₂ 及び k₃ は、状態コストに対する重みである。もし合流車両がランプ上で後続車両と先行車両との間である（すなわち、d x₁₂ < 0、及び d x₁₂ > d x₀₂ という条件にある）なら、k₁ = 1、k₂ = 10、及び k₃ = 10 であり、さもなければ、k₁ = 10、k₂ = 10、及び k₃ = 0 である。コストは、自動車は、後続車両と先行車両との中間に、後続車両と同一の速度で合流することを誘起すべく設計される。初期状態は、-100 < d x₁₂ < 100 メートル、-10 < d v₁₂ < 10 メートル / 秒、-100 < d x₀₂ < -5 メートル、及び -10 < d v₀₂ < 10 メートル / 秒において、ランダムに選択される。Z 値を近似するのに、ガウス放射基底関数が使用された。

【 数 2 1 】

$$f_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}_i(\mathbf{x} - \mathbf{m}_i)\right)$$

式中、 \mathbf{m}_i 及び \mathbf{S}_i は、第 i 番目の放射基底関数に対する平均及び逆共分散である。高速道路合流のシミュレーションに対し、 Z 値は、4,096 個のガウス放射基底関数であって、それらの平均が、状態の次元毎に 8 個の値から成る格子の頂点上に設定されたガウス放射基底関数により近似された。各基底の標準偏差は、各次元において最も近い 2 つの基底間の距離の 0.7 であった。上記例において

10

$\rho(\mathbf{x})$

の実際の値は一定であるため、

$\rho(\mathbf{x})$

を推定するのに $g(\mathbf{x}) = 1$ が使用された。上記方法は、受動的動力学的值をシミュレートすることにより収集された 10,000 個のサンプルから、ポリシーを最適化した。図 6 は、本明細書中に記述された方法により決定された順次的な制御入力を用い、125 個の異なる初期状態から開始し、(収束に対して必要とされる反復の回数として表現された) 30 秒以内に好首尾に合流する割合を示している。

20

【0069】

上記の詳細な説明においては、その一部を構成する添付図面に対する参照が為されている。図において、同様の記号は典型的に、状況が別様に示唆するのでなければ、同様の構成要素を特定している。詳細な説明、図、及び、請求項中に記述された代表的実施形態は、限定的であることを意味しない。本明細書中に呈示された主題の有効範囲から逸脱せずに、他の実施形態が利用され得ると共に、他の変更が為され得る。概略的に本明細書中に記述されると共に各図中に示された本開示の各見地は、全てが本明細書において明示的に企図された多様な異なる構成にて、配置、置換、結合、分離、及び、設計され得ることは容易に理解される。

30

【0070】

本開示を読んだに当業者により理解され得るように、本明細書中に記述された種々の見地は、方法、コンピュータシステム、又は、コンピュータプログラム製品として具現され得る。従って、それらの見地は、完全にハードウェアの実施形態、完全にソフトウェアの実施形態、又は、ソフトウェア及びハードウェアの見地を組み合わせた実施形態の形態を取り得る。更に、斯かる見地は、一種類以上のコンピュータ可読記憶媒体であって、本明細書中に記述された機能を実行するために当該記憶媒体内に又は当該記憶媒体上に具現されたコンピュータ可読プログラムコード又は命令を有するコンピュータ可読記憶媒体により記憶されたコンピュータプログラム製品の形態を取り得る。これに加え、本明細書中に記述されたデータ、命令又は事象を表す種々の信号は、送信元と送信先との間に、金属ワイヤ、光ファイバ、及び/又は、(例えば、空気及び/又は空間などの)無線送信媒体のような信号導通媒体を通して進行する電磁波の形態で伝達され得る。

40

【0071】

本明細書中で用いられるように、「一つの (a)」及び「一つの (an)」という語句は、一つ、又は、一つより多いものとして定義される。本明細書中で用いられるように、「複数の」という語句は、2 つ、又は、2 つより多いものとして定義される。本明細書中で用いられるように、「別の (another)」という語句は、少なくとも第 2 のもの、又は、そ

50

れより多いものとして定義される。本明細書中で用いられるように、「含む」及び／又は「有する」という語句は、備える（すなわち非制限的表現）として定義される。本明細書中で用いられるように、「～及び～の少なくとも一つ」という語句は、関連して列挙された対象物のうちの一つ以上の対象物の任意の全ての可能な組み合わせを参照かつ包含する。一例として、「A、B及びCの少なくとも一つ」という表現は、Aのみ、Bのみ、Cのみ、又は、（例えば、AB、AC、BC又はABCなどの）それらの任意の組合せを包含する。

【0072】

従って、本発明の有効範囲を表すものとしては、上述の明細書ではなく、以下の各請求項に対して参照が為されるべきである。

本明細書に開示される発明は以下の実施態様を含む。

（１）車両の自律的動作を適応的に制御するコンピュータ実行型方法であって、該方法は、

a) 車両を自律的に制御するように構成されたコンピュータ処理システムにおけるcriticネットワークにおいて、受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、actorネットワークにより適用されたときに車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、

b) コンピュータ処理システム内においてcriticネットワークに対して作用的に連結されたactorネットワークにおいて、車両に対して適用されて到達コストに対する最小値を生成する制御入力決定すること、とを備え、

前記actorネットワークは、推定平均コストと、近似された到達コスト関数から決定された推定到達コストと、車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成される、方法。

（２）前記近似された到達コスト関数は、以下の関係に従い、重み付けされた放射基底関数の線形結合を用いて決定され、

【数 2 2】

$$\hat{Z}(x) := \sum_{j=0}^N \omega_j f_j(x)$$

式中、 ω_j は重みであり、 f_j は第 j 番目の放射基底関数であり、 N は近似された到達コスト関数を決定するために使用される放射基底関数の個数であり、且つ、

$\hat{Z}(x)$

は近似された到達コスト関数である、上記（１）に記載の方法。

（３）前記近似された到達コスト関数において使用される重み ω_j は、以下の関係に従って更新され、

【数 2 3】

$$\omega^{i+1} = \tilde{\omega}^i + \lambda_1 + \sum_{l=0}^N \lambda_{2l} \delta_{jl} + \lambda_3 \hat{Z}_{avg} f_k$$

式中、 δ_{jl} はディラックのデルタ関数を表し、上付き文字 i は反復の回数を表し、 λ_1 、 λ_2 、 λ_3 はラグランジュ乗数であり、且つ

\hat{Z}_{avg}

は推定平均コストである、上記（２）に記載の方法。

（４）ベルマン方程式の線形化版を用いて決定された近似的な時間差的誤差を用いて前記criticネットワークのパラメータを更新する段階を更に備える、上記（１）に記載の方法。

（５）前記criticネットワークパラメータを更新する段階は、前記車両が運動しているときに実施される、上記（４）に記載の方法。

（６）前記criticネットワークにより決定された推定平均コストは、以下の関係に従って更新され、

【数２４】

$$\hat{Z}_{avg}^{i+1} = \hat{Z}_{avg}^i - \alpha_2^i e_k \hat{Z}_k$$

式中、 α_2^i は学習率であり、 e_k は近似的な時間差的誤差であり、

\hat{Z}_k

は前記近似された到達コスト関数から決定された推定コストであり、

\hat{Z}_{avg}^i

は状態 i における推定平均コストであり、且つ、

\hat{Z}_{avg}^{i+1}

は状態 $i + 1$ における推定平均コストである、上記（４）に記載の方法。

（７）前記受動的に収集されたデータは、前記criticネットワークパラメータを更新する間に使用される唯一のデータである、上記（４）に記載の方法。

（８）以下の関係に従って決定された近似的な時間差的誤差を用いて前記criticネットワークのパラメータを更新する段階を更に備え、

【数２５】

$$e_k := \hat{Z}_{avg} \hat{Z}_k - \exp(-q_k) \hat{Z}_{k+1}$$

式中、 e_k は近似的な時間差的誤差であり、

\hat{Z}_{avg}

は推定平均コストであり、

\hat{Z}_k

10

20

30

40

50

は状態 k における推定到達コストであり、

$$\hat{z}_{k+1}$$

は状態 $k + 1$ における推定到達コストであり、且つ、 q_k は状態 k における状態コストである、上記 (1) に記載の方法。

(9) 前記近似された到達コスト関数は、前記 critic ネットワークにおいてリアルタイムで学習される、上記 (1) に記載の方法。

(10) 前記ノイズレベルは、以下の関係に従い、重み付けされた基底関数の線形結合を用いて学習され、

【数 2 6】

$$\rho(x; \mu) \approx \hat{\rho}(x; \mu) := \sum_j^M \mu_j g_j(x)$$

式中、 μ_j は推定ノイズレベルであり、 μ_j は、 g_j により表された第 j 番目の放射基底関数に対する重みであり、且つ、 M は、ノイズレベルを推定するために使用されるべき放射基底関数の個数である、上記 (1) に記載の方法。

(11) 以下の関係に従って決定された近似誤差を用いて前記 actor ネットワークの重み付けパラメータを更新する段階を更に備え、

【数 2 7】

$$d_k \approx q_k \Delta t - \hat{V}_{k+1} + \hat{V}_k + \hat{V}_{avg} + L_{k,k+1} \rho_k,$$

式中、 d_k は近似誤差であり、 q_k は状態 k における状態コストであり、

$$\hat{V}_k$$

は状態 k において近似された到達コストであり、

$$\hat{V}_{k+1}$$

は状態 $k + 1$ において近似された到達コストであり、

$$\hat{V}_{avg}$$

は近似された平均コストであり、且つ、

【数 2 8】

$$L_{k,k+1} := (0.5 \hat{V}_k - \hat{V}_{k+1})^\top B_k B_k^\top \hat{V}_k \Delta t$$

であり、式中、 B_k は状態 k における制御用動力学的値である、上記 (10) に記載の方法

10

20

30

40

50

—
 (1 2) 前記actorネットワークの重み付けパラメータの更新は、前記車両が運動しているときに実施される、上記 (1 1) に記載の方法。
 (1 3) 前記actorネットワークの重み付けパラメータは、以下の関係に従って更新され

—
 【数 2 9】

$$\mu^{i+1} = \mu^i - \beta^i d_k L_{k,k+1} g_k,$$

10

式中、
 μ^{i+1}

は状態 $i + 1$ における重み付けパラメータの値であり、

μ^i

20

は状態 i における重み付けパラメータの値であり、 β^i は学習率であり、 d_k は時間差的誤差であり、 g は放射基底関数である、上記 (1 1) に記載の方法。

(1 4) 記制御入力を用い、前記自律的動作を制御すべく使用可能な制御ポリシーを修正する段階を更に備える、上記 (1) に記載の方法。

(1 5) 前記推定平均コストが収束するまで、前記段階 (a) 及び (b) を反復的に実施して前記制御入力を再決定することにより、前記自律的動作を制御するために使用可能な制御ポリシーを最適化する段階を更に備える、上記 (1) に記載の方法。

(1 6) 前記制御ポリシーは、能動的探索なしで最適化される、上記 (1 5) に記載の方法。

(1 7) 車両の自律的動作を適応的に制御するように構成されたコンピュータ処理システムであって、該コンピュータ処理システムは、該コンピュータ処理システムの動作を制御する一つ以上のプロセッサと、該一つ以上のプロセッサにより使用可能なデータ及びプログラム命令を記憶するメモリとを備え、

30

前記一つ以上のプロセッサは、前記メモリ内に記憶された命令を実行して、

a) 受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コストと、前記車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定し、且つ、

b) 前記車両に対して適用されて前記到達コストに対する最小値を生成する制御入力を決定する、ように構成され、

前記一つ以上のプロセッサは、前記推定平均コストと、前記近似された到達コスト関数から決定された到達コストと、前記車両の現在の状態に対する制御用動力学的値と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力を決定するように構成される、コンピュータ処理システム。

40

(1 8) 前記一つ以上のプロセッサは、前記メモリ内に記憶された命令を実行し、前記推定平均コストが収束するまで、前記段階 (a) 及び (b) を反復的に実施して前記制御入力を再決定することにより、前記自律的動作を制御するために使用可能な制御ポリシーを最適化するように構成される、上記 (1 7) に記載のコンピュータ処理システム。

(1 9) コンピュータシステムにより実行可能な命令が自身内に記憶された、一時的でないコンピュータ可読媒体であって、

a) 受動的に収集されたデータのサンプルと、状態コストとを用いて、推定平均コスト

50

と、車両の到達コストに対する最小値を生成する近似された到達コスト関数とを決定することと、

b) 前記車両に対して適用されて前記到達コストに対する最小値を生成する制御入力を決定すること、とを備える機能を実施させ、

前記制御入力は、前記到達コストに対する最小値を生成し、且つ前記平均コストと、前記近似された到達コスト関数から決定された到達コストと、前記車両の現在の状態に対する制御用動力学的值と、受動的に収集されたデータのサンプルとを用いて、ノイズレベルを推定することにより制御入力が決定される、一時的でないコンピュータ可読媒体。

(20) 前記命令は、前記推定平均コストが収束するまで、前記段階(a)及び(b)を反復的に繰り返して前記制御入力を再決定することにより、前記自律的動作を制御するために使用可能な制御ポリシーを最適化するように実行可能である、上記(19)に記載の一時的でないコンピュータ可読媒体。

【図面】

【図1】

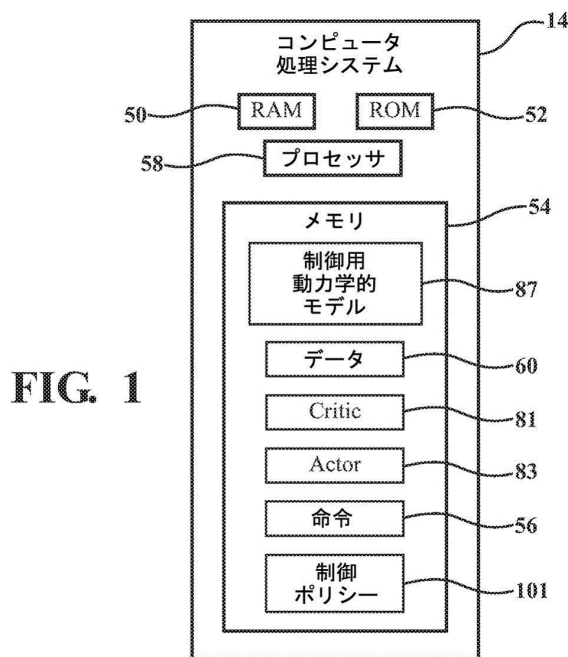


FIG. 1

【図2】

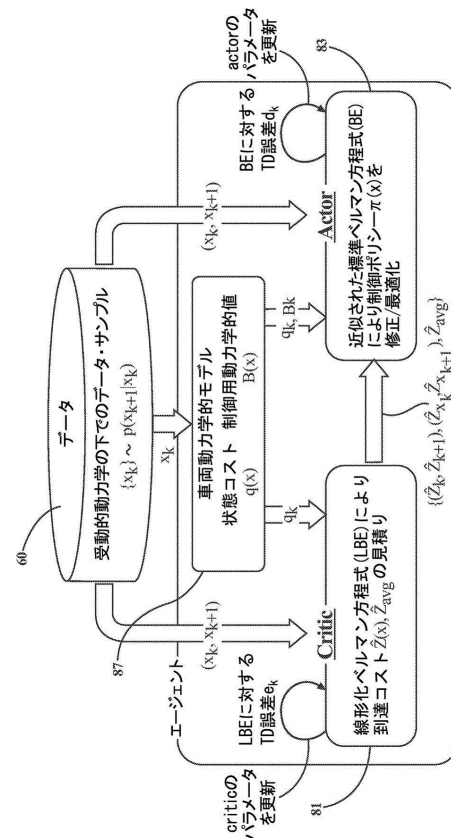


FIG. 2

10

20

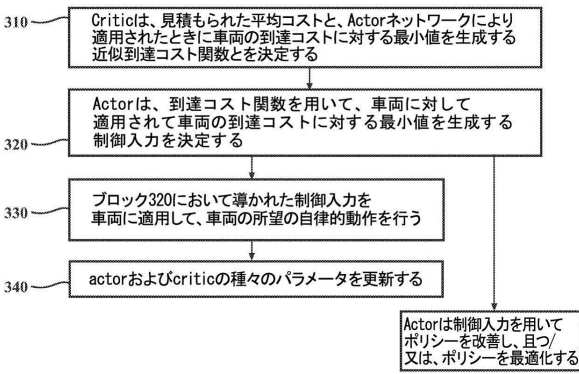
30

40

50

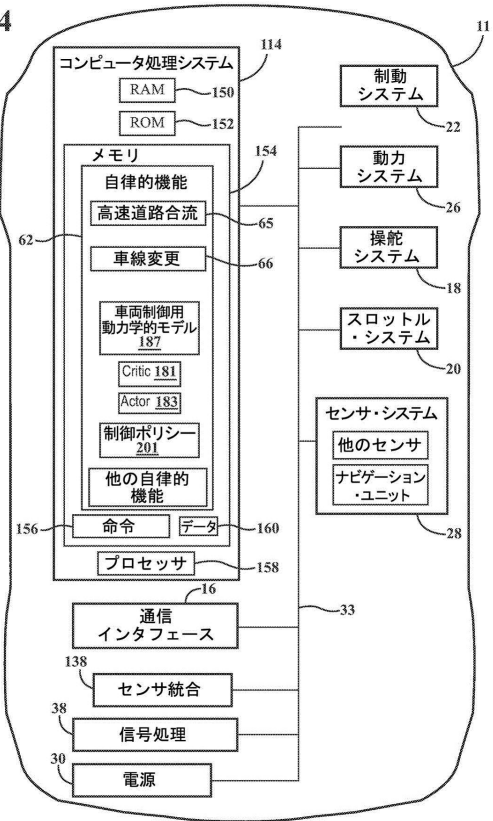
【図 3】

FIG. 3



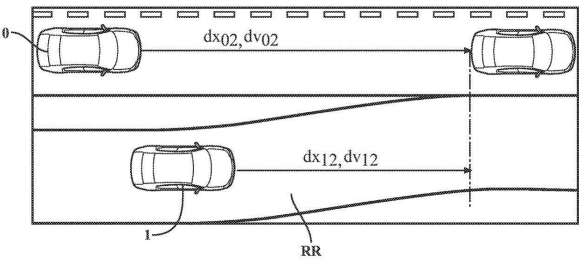
【図 4】

FIG. 4



【図 5】

FIG. 5



【図 6】

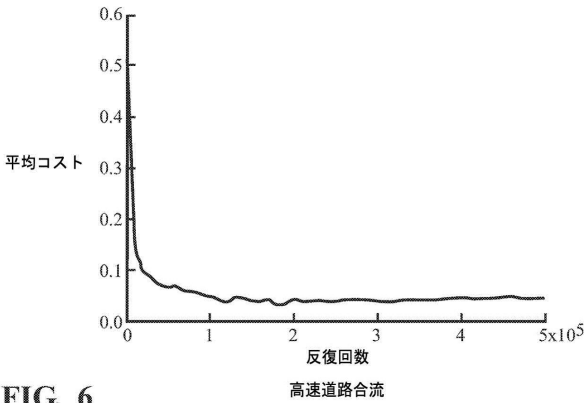


FIG. 6

フロントページの続き

(74)代理人 100123593

弁理士 関根 宣夫

(72)発明者 西 智樹

アメリカ合衆国, ケンタッキー 41018, アーランガー, アトランティック アベニュー 25,
シーノールトヨタ モーター エンジニアリング アンド マニュファクチャリング ノース アメリ
カ, インコーポレイティド

審査官 杉山 悟史

(56)参考文献 特開平10-254505(JP, A)

特開2006-313512(JP, A)

米国特許出願公開第2013/0262353(US, A1)

米国特許出願公開第2016/0092764(US, A1)

(58)調査した分野 (Int.Cl., DB名)

G05B 1/00 - 7/04

11/00 - 13/04

17/00 - 17/02

21/00 - 21/02