

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4644833号
(P4644833)

(45) 発行日 平成23年3月9日(2011.3.9)

(24) 登録日 平成22年12月17日(2010.12.17)

(51) Int.Cl. F I
B 2 5 J 5/00 (2006.01) B 2 5 J 5/00 E
 B 2 5 J 5/00 F

請求項の数 3 (全 34 頁)

(21) 出願番号	特願2004-105291 (P2004-105291)	(73) 特許権者	393031586
(22) 出願日	平成16年3月31日 (2004. 3. 31)		株式会社国際電気通信基礎技術研究所
(65) 公開番号	特開2005-288594 (P2005-288594A)		京都府相楽郡精華町光台二丁目2番地2
(43) 公開日	平成17年10月20日 (2005.10.20)	(74) 代理人	110001195
審査請求日	平成19年3月14日 (2007. 3. 14)		特許業務法人深見特許事務所
		(74) 代理人	100064746
			弁理士 深見 久郎
		(74) 代理人	100085132
			弁理士 森田 俊雄
		(74) 代理人	100083703
			弁理士 仲村 義平
		(74) 代理人	100096781
			弁理士 堀井 豊
		(74) 代理人	100098316
			弁理士 野田 久登

最終頁に続く

(54) 【発明の名称】 2足歩行移動装置

(57) 【特許請求の範囲】

【請求項 1】

右上腿部と前記右上腿部に右膝の関節を介して接続する右下腿部とを備える右脚と、
 左上腿部と前記左上腿部に左膝の関節を介して接続する左下腿部とを備える左脚と、
 前記右脚または前記左脚の接地を検出し、前記右上腿部に対して前記右下腿部が成す右膝角と前記左上腿部に対して前記左下腿部が成す左膝角とを検出するための第1のセンサ群とを備え、

前記右脚および前記左脚に接続し、前記右脚および前記左脚を駆動して2足歩行を行わせるための腰部と、

前記腰部のピッチ角を検出するためのピッチ角センサと、

前記腰部と前記右上腿部との成す右関節角と前記腰部と前記左上腿部との成す左関節角とを検出するための第2のセンサ群とをさらに備え、

前記腰部は、

前記右下腿部および前記左下腿部に対する状態マシンに基づいて、前記右下腿部および前記左下腿部を駆動する下腿制御信号を生成する下腿制御部と、

前記第2のセンサ群と前記ピッチ角センサの検知結果を受けて、前記右上腿部および前記左上腿部を駆動するための上腿制御信号を生成するための動的制御装置とを備え、

前記動的制御装置は、

前記右上腿部および前記左上腿部にそれぞれ対応する周期的な時間発展を行う内部状態を有するセントラルパターンジェネレータに対して、強化学習に基づいて前記右上腿部お

よび前記左上腿部の状態推定を行い、前記内部状態に対応する目標値に対するPDサーボ系の出力として、前記右上腿部および前記左上腿部に対する前記上腿制御信号を生成する制御手段を含み、

前記上腿制御信号および前記下腿制御信号に基づいて、前記右脚および前記左脚を駆動するための駆動手段をさらに備え、

前記下腿制御部の前記状態マシンは、

前記右脚および前記左脚のそれぞれについて、第1の膝屈曲状態と、第1の膝伸長状態と、第2の膝屈曲状態と、第2の膝伸長状態とを順次遷移する状態マシンであって、

前記第1の膝屈曲状態から前記第1の膝伸長状態への遷移および前記第2の膝屈曲状態から前記第2の膝伸長状態への遷移は、前記第2のセンサ群の検知結果に基づいて、前記左上腿部と前記右上腿部の成す角度が所定値を下回ることを条件として遷移が発生し、

前記第1の膝伸長状態から前記第2の膝屈曲状態への遷移および前記第2の膝伸長状態から前記第1の膝屈曲状態への遷移は、前記第1のセンサ群の検知結果に基づいて、前記左脚または前記右脚の接地に応じて発生し、

前記駆動手段は、前記状態マシンの状態に応じて予め定められた前記右膝角および前記左膝角についての目標角度と前記第1のセンサ群により検知される前記右膝角および前記左膝角とに基づくPDサーボ系の出力として、前記右膝および前記左膝とに加えるトルクを出力する膝駆動手段を含む、2足歩行移動装置。

【請求項2】

前記制御手段は、

前記右上腿部および前記左上腿部の状態情報に基づいて、価値関数と前記強化学習中に得られる報酬系列とに基づいて、フィードバックパラメータを更新するフィードバック制御器を含み、

前記フィードバック制御器からの出力に応じて変化する前記セントラルパターンジェネレータの前記内部状態に基づいて、前記上腿制御信号を生成する、請求項1記載の2足歩行移動装置。

【請求項3】

前記強化学習は、方策勾配法により行われる、請求項2記載の2足歩行移動装置。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、非線形の制御対象である2足歩行の周期運動を実現する上で有効となる動的制御装置を用いた2足歩行移動装置に関する。

【背景技術】

【0002】

ロボットなどを制御しようとするとき、センサノイズや、そもそもセンサを配置することができないことによって、制御のために必要な状態変数を直接的には測定できない状況が考えられる。

【0003】

そのような場合、状態観測器（オブザーバ）（たとえば、非特許文献1を参照）やカルマンフィルタ（たとえば、非特許文献2、非特許文献3を参照）を用いることが一般的である。しかしながら、対象のダイナミクスが非線形である場合、これらの手法では隠れ状態の推定は困難な場合がある。近年、非線形系に適用可能なオブザーバの提案がなされているものの、それぞれ特定の条件を満たさなければ適用できないなどの問題がある（たとえば、非特許文献4、非特許文献5を参照）。

【0004】

また、拡張カルマンフィルタ（局所線形モデルを用いて状態分布の更新を行う）（たとえば、非特許文献3を参照）やモンテカルロフィルタ（モンテカルロ法により生成した多数の粒子により状態分布を近似する）（たとえば、非特許文献6を参照）などの手法が非線形系での状態推定法として知られているが、いずれも状態分布を陽に求めなければならな

10

20

30

40

50

い。

【 0 0 0 5 】

すなわち、従来の制御器は、基本的には、現在の状態を観測し、そこからの直接の写像によって制御出力を与える必要がある。

【非特許文献 1】D.G. Luenberger著、“An introduction to observers”、IEEE Trans., AC, Vol.16, pp. 596--602, 1971.

【非特許文献 2】R.E. Kalman and R.S. Bucy著、“New results in linear filtering and prediction theory”、Trans., ASME, Series D, J. of Basic Engineering, Vol.83, No.1, pp. 95--108, 1961.

【非特許文献 3】F.L. Lewis著、“Optimal Estimation: with an Introduction to Stochastic Control Theory”、John Wiley & Sons, 1977. 10

【非特許文献 4】志水清孝, 鈴木俊輔, 田中哲史著、“こう配降下法による非線形オブザーバ(非線形システムの状態観測器)”、電子情報通信学会論文誌 A, Vol. J83-A, No.8, pp. 956--964, 2000.

【非特許文献 5】H.Nijmeijer and T.L. Fossen著、“Directions in Nonlinear Observer Design”、Springer-Verlag, London, 1999.

【非特許文献 6】G.Kitagawa著、“Monte Carlo Filter and Smoother for Non-Gaussian Nonlinear State Models”、Journal of Computational and Graphical Statistics, Vol.5, pp. 1--25, 1996.

【発明の開示】 20

【発明が解決しようとする課題】

【 0 0 0 6 】

本発明の目的は、制御器に内部状態とダイナミクスを持たせることで、周期運動に対し状態観測器の学習を容易とすることが可能な動的制御装置を用いた 2 足歩行移動装置を提供することである。

【課題を解決するための手段】

【 0 0 0 7 】

本発明では、学習の目的は、ある瞬間での推定誤差を少なくするのではなく、タスクを行っている期間を通じての推定誤差を少なくすることである。そこで、方策勾配法(強化学習)の枠組みを用いて非線形状態観測器の構築を行っている。 30

【 0 0 0 8 】

本発明の 1 つの局面では、状態観測器のダイナミクスへの入力を学習器の行動と考え、現在の観測可能な状態(観測出力)、状態観測器の状態、制御器の出力から、適切に状態観測器の状態を制御対象の状態へと導く行動則を方策勾配法によって獲得している。

【 0 0 0 9 】

したがって、この発明の 1 つの局面に従うと、2 足歩行移動装置であって、右上腿部と右上腿部に右膝の関節を介して接続する右下腿部とを備える右脚と、左上腿部と左上腿部に左膝の関節を介して接続する左下腿部とを備える左脚と、右脚または左脚の接地を検出し、右上腿部に対して右下腿部が成す右膝角と左上腿部に対して左下腿部が成す左膝角とを検出するための第 1 のセンサ群と、右脚および左脚に接続し、右脚および左脚を駆動して 2 足歩行を行わせるための腰部と、腰部のピッチ角を検出するためのピッチ角センサと、腰部と右上腿部との成す右関節角と腰部と左上腿部との成す左関節角とを検出するための第 2 のセンサ群とをさらに備え、腰部は、右下腿部および左下腿部に対する状態マシンに基づいて、右下腿部および左下腿部を駆動する下腿制御信号を生成する下腿制御部と、第 2 のセンサ群とピッチ角センサの検知結果を受けて、右上腿部および左上腿部を駆動するための上腿制御信号を生成するための動的制御装置とを備え、前記動的制御装置は、前記右上腿部および前記左上腿部にそれぞれ対応する周期的な時間発展を行う内部状態を有するセントラルパターンジェネレータに対して、強化学習に基づいて右上腿部および左上腿部の状態推定を行い、内部状態に対応する目標値に対する P D サーボ系の出力として、右上腿部および左上腿部に対する上腿制御信号を生成する制御手段を含み、上腿制御信号 40 50

および下腿制御信号に基づいて、右脚および左脚を駆動するための駆動手段をさらに備え、下腿制御部の状態マシンは、右脚および左脚のそれぞれについて、第1の膝屈曲状態と、第1の膝伸長状態と、第2の膝屈曲状態と、第2の膝伸長状態とを順次遷移する状態マシンであって、第1の膝屈曲状態から第1の膝伸長状態への遷移および第2の膝屈曲状態から第2の膝伸長状態への遷移は、第2のセンサ群の検知結果に基づいて、左上腿部と右上腿部の成す角度が所定値を下回ることを条件として遷移が発生し、第1の膝伸長状態から第2の膝屈曲状態への遷移および第2の膝伸長状態から第1の膝屈曲状態への遷移は、第1のセンサ群の検知結果に基づいて、左脚または右脚の接地に応じて発生し、駆動手段は、状態マシンの状態に応じて予め定められた右膝角および左膝角についての目標角度と第1のセンサ群により検知される右膝角および左膝角とに基づくPDサーボ系の出力として、右膝および左膝とに加えるトルクを出力する膝駆動手段を含む。

10

【0010】

好ましくは、制御手段は、右上腿部および左上腿部の状態情報に基づいて、価値関数と強化学習中に得られる報酬系列とに基づいて、フィードバックパラメータを更新するフィードバック制御器を含み、フィードバック制御器からの出力に応じて変化するセントラルパターンジェネレータの内部状態に基づいて、上腿制御信号を生成する。

【0011】

好ましくは、強化学習は、方策勾配法により行われる。

【発明の効果】

【0015】

本発明では、行動則自体が内部変数とその微分方程式によって表されるダイナミクスを持つため、行動則と物理系の引き込みの性質を利用することができ、周期運動に対して状態推定を行う場合に、状態観測器の学習を容易とすることが可能である。さらに、センサ入力のノイズや時間遅れに対しても、ロバストな性質を持つ制御器を実現することができる。

20

【0016】

さらに、本発明によれば、2足歩行移動装置の自由度が増加し、状態空間が高次元となった場合でも、周期運動に対し状態観測器の学習を容易とすることが可能な動的制御装置動的制御装置を用いた2足歩行移動装置を提供することができる。

【発明を実施するための最良の形態】

30

【0017】

以下、図面を参照して本発明の実施の形態について説明する。

【0018】

以下の説明では、一例として、本発明を2足歩行の制御に適用する場合を説明するが、本発明は、必ずしもこのような場合に限定されるものではなく、たとえば、より一般的に周期運動を行う系に対して有効な制御システムを提供するものである。特に、本発明は、周期運動を行う劣駆動（機械）系に適用するのに適した制御システムを提供する。

【0019】

[実施の形態1]

(本発明のシステム構成)

図1は、本発明の動的制御装置を用いた2足歩行移動システム1000の一例を示す概念図である。

40

【0020】

図1を参照して、システム1000は、動的制御装置100と、動的制御装置100の上部に設けられる胴部40と、動的制御装置100により駆動制御される脚部とを備える。脚部は、右脚10rと左脚10lとを有し、各脚は、接地面近傍に設けられるセンサ20rおよび20lとを備える。一方、胴部40には、センサ30が設けられる。

【0021】

センサ20rは、胴部40の中心線4に対する右脚の角度 θ_r 、角速度 $d\theta_r/dt$ という情報を検出し、また、センサ20lは、中心線4に対する左脚の角度 θ_l 、角速度 $d\theta_l/dt$

50

/d t という情報を検出し、それぞれ、動的制御装置 100 に通知する。さらに、センサ 30 は、鉛直方向 2 に対する胴部 40 のピッチ角 ρ 、角速度 $d\rho/dt$ という情報を検出し、それぞれ、動的制御装置 100 に通知する。

【0022】

動的制御装置 100 は、センサ 20r、センサ 20l、センサ 30 からの情報に基づいて、右脚 10r および左脚 10l の動作を制御する。

【0023】

図 2 は、図 1 に示した動的制御装置 100 の構成を示すブロック図である。

【0024】

図 2 を参照して、動的制御装置 100 は、センサ 20r、センサ 20l、センサ 30 からの信号を受け取る通信インタフェース 106 と、後に説明する制御パラメータやセンサからの情報を格納しておくための記憶装置 104 と、センサ 20r、センサ 20l、センサ 30 からの情報を用いて学習して獲得した動的行動則に基づき、制御信号を生成する演算処理部 102 と、演算処理部 102 からの制御信号に基づいて、右脚 10r および左脚 10l の駆動制御を行うための駆動部 108 とを備える。

10

【0025】

以下では、動的制御装置 100 の制御動作のための準備の処理および制御動作について説明する。

【0026】

(1-1. 動的行動則)

20

まず、本発明の制御動作を説明する前提として、「動的行動則」について説明する。

【0027】

図 3 は、このような動的行動則を説明するための概念図である。

【0028】

「動的行動則」とは、図 3 (1) に示すような行動則自体が内部変数とその微分方程式によって表されるダイナミクスを持つ枠組である。

【0029】

これをより具体的に表現すると、図 3 (2) に示すように観測器が内部変数及びその微分方程式を持つようなものや、図 3 (3) に示すような制御器が内部変数及びその微分方程式を持つようなものである場合が考えられる。

30

【0030】

このような「動的行動則」に基づいて、制御対象を制御するような制御装置を「動的制御装置」と呼ぶことにする。

【0031】

以下では、まず、図 3 (3) の枠組を用いることで、動的制御装置 100 により、2 足歩行運動を実現する場合を考える。

【0032】

行動則が内部状態を持つことによって、行動則と物理系の引き込みの性質を利用することが出来るため、周期運動や状態推定を行う場合に有効であると考えられる。さらに、センサ入力のノイズや時間遅れに対しても、ある程度ロバストな性質を持つことが期待できる。

40

【0033】

このような行動則を学習によって獲得する場合、行動則の内部状態が隠れ変数となることは問題となる。しかし、行動則の内部状態が物理系の状態に対して引き込むことで、それぞれの状態は一意に対応するようになる。ただし、過渡的な状態では隠れ状態を扱う必要がある。

【0034】

そこで本発明では、動的行動則の獲得手法として、後に説明するように、隠れ変数が存在する環境においても適用可能な方策勾配法を用いる。

【0035】

50

以下では、本発明の学習システムおよび、動的行動則を構成するセントラルパターンジェネレータ (central pattern generator : C P G) と C P G へのフィードバック制御器について説明を行う。

【 0 0 3 6 】

(1 - 2 . 学習システム)

以下の説明では、周期運動の例として、3リンク2足歩行ロボットモデルを用いた2足歩行運動に対して動的行動則を適用する。

【 0 0 3 7 】

図4は、図2で説明した演算処理部102の行う処理を示す機能ブロック図である。以下に説明するとおり、演算処理部102は、学習システムとして機能する。

10

【 0 0 3 8 】

図4に示すとおり、この学習システムは、基本的に、C P G 処理部1026とフィードバック制御器1022によって動的行動則を構成する。学習に用いる状態 x は、以下の式で表される。

【 0 0 3 9 】

【 数 1 】

$$x = (x_1, x_2, x_3, x_4)^T = (\theta_l + \theta_p, \theta_r + \theta_p, \dot{\theta}_l + \dot{\theta}_p, \dot{\theta}_r + \dot{\theta}_p)^T$$

【 0 0 4 0 】

ただし、上述のとおり、 θ_r 、 θ_l は、それぞれロボットの鉛直方向からの左右の脚10r、10lの角度であり、 θ_p は胴体40のピッチ角である。

20

【 0 0 4 1 】

つまり、学習システムはロボットから直接得られる信号のみで状態空間を構成しており、C P G の内部状態を用いていないという特徴を有している。

【 0 0 4 2 】

また、ここでは全状態観測を仮定し、 $y = x$ であるとする。

【 0 0 4 3 】

(1 - 3 . セントラルパターンジェネレータ (C P G))

演算処理部102により実現される学習システムで、動的行動則を構成するC P G 処理部1026の構成として、以下の式で表される神経振動子モデルを用いる。なお、このような神経振動子モデルについては、たとえば、文献：Kiyoshi Matsuoka著、“Sustained oscillations generated by mutually inhibiting neurons with adaptation.”、Biological Cybernetics, Vol.52, pp. 367-376, 1985に開示がある。

30

【 0 0 4 4 】

【数2】

$$\tau \dot{z}_1 = -z_1 - \omega_{12}q_2 - \omega_{13}q_3 - \beta p_1 + z_0 + a_1 \quad (1)$$

$$\tau' \dot{p}_1 = -p_1 + q_1 \quad (2)$$

$$q_1 = \max(0, z_1) \quad (3)$$

$$\tau \dot{z}_2 = -z_2 - \omega_{21}q_1 - \omega_{24}q_4 - \beta p_2 + z_0 + a_2 \quad (4)$$

$$\tau' \dot{p}_2 = -p_2 + q_2 \quad (5)$$

$$q_2 = \max(0, z_2) \quad (6)$$

$$\tau \dot{z}_3 = -z_3 - \omega_{34}q_4 - \omega_{31}q_1 - \beta p_3 + z_0 + a_3 \quad (7)$$

$$\tau' \dot{p}_3 = -p_3 + q_3 \quad (8)$$

$$q_3 = \max(0, z_3) \quad (9)$$

$$\tau \dot{z}_4 = -z_4 - \omega_{43}q_3 - \omega_{42}q_2 - \beta p_4 + z_0 + a_4 \quad (10)$$

$$\tau' \dot{p}_4 = -p_4 + q_4 \quad (11)$$

$$q_4 = \max(0, z_4) \quad (12)$$

【0045】

ここで、変数： z 、 p はニューロン内の状態、 q はニューロンの出力、 z_0 は持続入力、定数はニューロンの疲労係数、 τ 、 τ' は、 z 、 p の時定数、 ω は拮抗ニューロン間の結合係数である。また、 a は、後に説明するフィードバック制御器からの出力項である。

【0046】

図5は、式(1)～(12)で表される神経振動子モデルによるCPGを示す概念図である。図5においては、ニューロン内の状態 z 、 p の間で、相互に正の結合を行うものは白丸で、負の結合を行うものは黒丸で示している。

【0047】

図6は、このようなCPGの出力 q を構成する変数 z_1 、 z_2 の波形を示す図である。なお、この計算では、例として、 $\omega_{12} = 0.05$ 、 $\omega_{21} = 0.6$ 、 $\omega_{13} = 2.5$ 、 $\omega_{31} = 2.0$ 、 $z_0 = 0.1$ 、 $a = 0$ を用いた。

【0048】

式(1)～(12)で表されるモデルにしたがって、CPGの内部状態変数 z_1 、 z_2 が周期的に変化していることがわかる。

【0049】

さらに、図4の学習システムのPDサーボ処理部1028では、以下に示すとおり、各ニューロンの出力の差を両脚のサーボ系の目標関節角 θ^d とした。

【0050】

【数3】

$$\theta_l^d = (q_1 - q_2) \quad (13)$$

$$\theta_r^d = (q_3 - q_4) \quad (14)$$

10

20

30

40

50

【 0 0 5 1 】

ただし、 θ_l^d は左脚の目標関節角、 θ_r^d は右脚の目標関節角である。

【 0 0 5 2 】

P Dサーボ処理部 1 0 2 8の結果出力されるロボットへのトルク入力 u は、次に示す P Dサーボ系の出力を用いる。

【 0 0 5 3 】

【数 4】

$$u_l = K_p(\theta_l^d - \theta_l) - K_d\dot{\theta}_l \quad (15)$$

$$u_r = K_p(\theta_r^d - \theta_r) - K_d\dot{\theta}_r \quad (16)$$

10

【 0 0 5 4 】

ただし、 u_l は左脚に対するトルク入力、 u_r は右脚に対するトルク入力である。また、 K_p は位置ゲイン、 K_d は速度ゲインである。

【 0 0 5 5 】

(1 - 4 . フィードバック制御器 1 0 2 2)

上述の C P G へのフィードバック制御器 1 0 2 2 は、次の確率分布 (1 7) によって表される。

【 0 0 5 6 】

20

【数 5】

$$\pi(v_j, \mathbf{x}; \mathbf{w}) = \frac{1}{\sigma_j \sqrt{2\pi}} \exp\left(-\frac{(v_j - \mu_j)^2}{2\sigma_j^2}\right) \quad (17)$$

【 0 0 5 7 】

ただし、 \mathbf{x} は制御対象の状態ベクトル、 \mathbf{w} はパラメータベクトルである。従って j 番目の出力の実現値 v_j は、以下の式 (1 8) によって与えられる。

【 0 0 5 8 】

【数 6】

30

$$v_j(t) = \mu_j(t) + \sigma_j n_j(t) \quad (18)$$

【 0 0 5 9 】

ただし、 $n_j(t) \sim N(0, 1)$ であり、 $N(0, 1)$ は平均 0、分散 1 の正規分布を表す。

【 0 0 6 0 】

ここでは出力を飽和させるために、出力飽和処理部 1 0 2 4 において、関数 $d(\cdot)$ を用いて以下の式 (1 9) のように、最終的な制御器の出力 $a_j(t)$ を決定する。

【 0 0 6 1 】

40

【数 7】

$$a_j(t) = a_j^{max} d(v_j(t)) \quad (19)$$

【 0 0 6 2 】

ただし、以下の説明では、一例として、 $d(\cdot)$ としては、以下の式を用いる。

【 0 0 6 3 】

【数 8】

$$d(x) = \frac{2}{\pi} \arctan\left(\frac{\pi}{2} x\right)$$

【0064】

ここでの a_j ($j=1 \sim 4$) は、式 (1) の左右の脚の神経振動子の伸筋、屈筋にそれぞれ対応する。

【0065】

(2. 方策勾配法)

「方策勾配法」とは、パラメータ化された確率的方策に従って行動選択を行い、方策を改善する方向に方策のパラメータを少しずつ更新する強化学習手法の 1 種である。以下に方策勾配法を用いた行動則の学習方法について述べる。

【0066】

(2-1. 連続時間・状態系でのテンポラル・ディファレンス (Temporal Difference) 誤差)

連続時間・状態系のダイナミクスを以下の式 (20) で表す。

【0067】

【数 9】

$$\frac{d\mathbf{x}(t)}{dt} = f(\mathbf{x}(t), \mathbf{u}(t)) \quad (20)$$

【0068】

ただし、 $\mathbf{x} \in \mathbb{R}^n$ は状態、 $\mathbf{u} \in \mathbb{R}^m$ は制御入力を表す。

【0069】

報酬は状態と制御入力の関数として、以下の式 (21) で与えられるとする。

【0070】

【数 10】

$$r(t) = r(\mathbf{x}(t), \mathbf{u}(t)) \quad (21)$$

【0071】

ある制御則 $(\mathbf{u}(t) | \mathbf{x}(t))$ のもとで、状態 $\mathbf{x}(t)$ の価値関数を以下の式 (22) で定義する。

【0072】

【数 11】

$$V^\pi(\mathbf{x}(t)) = E \left\{ \int_t^\infty e^{-\frac{s-t}{\tau}} r(\mathbf{x}(s), \mathbf{u}(s)) ds \mid \pi \right\} \quad (22)$$

【0073】

ただし、 τ は価値関数の時定数である。また、式 (22) の両辺の時間微分から、以下の式 (23) という拘束条件が与えられる。

【0074】

【数 12】

$$\frac{dV^\pi(\mathbf{x}(t))}{dt} = \frac{1}{\tau} V^\pi(\mathbf{x}(t)) - r(t) \quad (23)$$

10

20

30

40

50

【 0 0 7 5 】

$V(x(t)) = V(x(t); w)$ を価値関数の予測値とする。ただし、 w は評価値の予測値のパラメータである。

【 0 0 7 6 】

予測が正しければ、式 (2 3) を満たす。予測が正しくない場合、下式 (2 4) に示した予測誤差を減らすように学習を行う。

【 0 0 7 7 】

【数 1 3】

$$\delta(t) = r(t) - \frac{1}{\tau} V(t) + \dot{V}(t) \quad (24)$$

10

【 0 0 7 8 】

上式は連続時間系での TD 誤差である。

【 0 0 7 9 】

(2 - 2 . 方策勾配法の一般論)

動的計画法やグリーディ方策 (greedy policy) などの価値関数の評価を基に学習を行う場合では、環境がマルコフ決定過程である必要があるが、実問題に適用する場合には、ノイズやセンサの能力によってマルコフ決定過程を保証することは困難である。しかし、方策勾配法は、価値関数と共に、試行中に得られた累積報酬系列を考慮することで、環境が非マルコフ決定過程 (POMDP) でも適用することが出来る。

20

【 0 0 8 0 】

ここで、パラメータ w を持つ方策 π_w を用いた場合、以下の式 (2 5) が成り立つ。

【 0 0 8 1 】

【数 1 4】

$$\frac{\partial}{\partial w_i} E \{ V(t) | \pi_w \} = E \{ \delta(t) e_i(t) \} \quad (25)$$

【 0 0 8 2 】

ただし、以下の式 (2 6) が成り立つ。

30

【 0 0 8 3 】

【数 1 5】

$$\dot{e}_i(t) = -\frac{1}{\kappa} e_i(t) + \frac{\partial \ln \pi_w}{\partial w_i} \quad (26)$$

【 0 0 8 4 】

ここで、 κ はエリジビリティ・トレース (eligibility trace) の時定数である。テンポラル・ディファレンス誤差 $\delta(t)$ と方策のエリジビリティ・トレース $e(t)$ により、価値関数の方策パラメータ w に関する勾配の不偏推定量を求めることが出来ることが与えられている。このような方策勾配法については、たとえば、文献：木村元、小林重信、" Actor に適正度の履歴を用いた actor-critic アルゴリズム - 不完全な value-function のもとの強化学習 "、人工知能学会誌、Vol.15, No.2, pp. 267-275, 2000 に記載がある。

40

【 0 0 8 5 】

よって、パラメータの更新則は次の式 (2 7) のようになる。

【 0 0 8 6 】

【数 1 6】

$$\dot{w}_i = \eta \delta(t) e(t) \quad (27)$$

【0 0 8 7】

ただし、 η は学習率である。

【0 0 8 8】

(2 - 3 . 動的行動則の学習)

以下では、上述の方策勾配法を用いて、動的行動則の獲得を行う。

【0 0 8 9】

ここでは、式(17)に示したフィードバック制御器の学習を行うことで望みの動的行動則を獲得することを考える。

【0 0 9 0】

(2 - 3 - 1 . 価値関数の更新)

まず、価値関数処理部 1 0 3 2 において演算される、連続状態における価値関数の表現方法として、以下の式(28)による正規化ガウス関数ネットワーク(normalized Gaussian network: NGnet)を用いる。なお、正規化ガウス関数ネットワークについては、後に説明する。

【0 0 9 1】

【数 1 7】

$$V(\mathbf{x}(t)) = \sum_i w_i^c b_i^c(\mathbf{x}(t)) \quad (28)$$

【0 0 9 2】

ただし、 $b_i^c(\cdot)$ は、正規化処理部 1 0 3 0 において \mathbf{x} に施される基底関数であり、 w_i^c は価値関数のパラメータである。

【0 0 9 3】

パラメータ w_i^c に対するエリジビリティ・トレース e_i^c と、TD 誤差を用いたパラメータ w_i^c の更新式は、それぞれ以下の式(29)および(30)のようになる。

【0 0 9 4】

【数 1 8】

$$\dot{e}_i^c(t) = -\frac{1}{\kappa^c} e_i^c(t) + b_i^c(\mathbf{x}(t)) \quad (29)$$

$$\dot{w}_i^c(t) = \alpha \delta(t) e_i^c(t) \quad (30)$$

【0 0 9 5】

ただし、 α は価値関数の学習率、 κ^c はエリジビリティ・トレースの時定数である。

【0 0 9 6】

(2 - 3 - 2 . フィードバック制御器の更新)

式(17)に示した確率的なフィードバック制御器 1 0 2 2 を用いる場合、その j 番目の出力の平均 μ_j と標準偏差 σ_j に関するエリジビリティは式(26)右辺第2項と同様、それぞれ以下のように与えられる。

【0 0 9 7】

10

20

30

40

【数 19】

$$\frac{\partial \ln \pi}{\partial \mu_j} = \frac{a_t - \mu_j}{\sigma_j^2} \quad (31)$$

$$\frac{\partial \ln \pi}{\partial \sigma_j} = \frac{(a_t - \mu_j)^2 - \sigma_j^2}{\sigma_j^3} \quad (32)$$

【0098】

10

ここではさらに、以下の式(33)および(34)のように、平均 μ を正規化ガウス関数ネットワークによって表し、標準偏差 σ をシグモイド関数によって表す。

【0099】

【数 20】

$$\mu_j = \sum_i w_{ij}^\mu b_i^\mu(\mathbf{x}(t)) \quad (33)$$

$$\sigma_j = \frac{1}{1 + \exp(-w_j^\sigma)} \quad (34)$$

20

【0100】

ただし、ノーターションとしては以下のとおりである。

【0101】

【数 21】

b_i^μ : 基底関数、

w_{ij}^μ, w_j^σ : 式(17)に示したフィードバック制御器のパラメータ

30

【0102】

これらのパラメータに対応するエリジビリティは以下の式(35)および(36)のように求められる。

【0103】

【数 22】

$$\frac{\partial \ln \pi}{\partial w_{ij}^\mu} = \frac{\partial \ln \pi}{\partial \mu_j} \frac{\partial \mu_j}{\partial w_{ij}^\mu} = \frac{(a_t - \mu_j) b_i^\mu(\mathbf{x}(t))}{\sigma_j^2} \quad (35)$$

$$\frac{\partial \ln \pi}{\partial w_j^\sigma} = \frac{\partial \ln \pi}{\partial \sigma_j} \frac{\partial \sigma_j}{\partial w_j^\sigma} = \frac{((a_t - \mu_j)^2 - \sigma_j^2)(1 - \sigma_j)}{\sigma_j^2} \quad (36)$$

40

【0104】

上式(35)(36)と式(26)(27)を考慮すると、以下の式(37)および(38)のようなフィードバックパラメータの更新則が得られる。

【0105】

【数 2 3】

$$\dot{w}_{ij}^{\mu} = \beta^{\mu} \delta(t) e_{ij}^{\mu}(t) \quad (37)$$

$$\dot{w}_j^{\sigma} = \beta^{\sigma} \delta(t) e_j^{\sigma}(t) \quad (38)$$

【0 1 0 6】

ただし、ノーターションとしては以下のとおりである。

10

【0 1 0 7】

【数 2 4】

 $\beta^{\mu}, \beta^{\sigma}$: 学習率 $e_{ij}^{\mu}(t), e_j^{\sigma}(t)$: それぞれのパラメータのエリジビリティ・トレース

【0 1 0 8】

また、式(35)(36)において、パラメータ w_{ij} が分母となっていることにより、 w_{ij} が 0 へと近付くとエリジビリティが発散することが問題となる。そこでエリジビリティ・トレースの更新には式(26)の代わりに次式を用いる。

20

【0 1 0 9】

【数 2 5】

$$\dot{e}_{ij}^{\mu}(t) = -\frac{1}{\kappa^{\mu}} e_{ij}^{\mu}(t) + \sigma_j^2 \frac{\partial \pi_w}{\partial w_{ij}} \quad (39)$$

$$\dot{e}_j^{\sigma}(t) = -\frac{1}{\kappa^{\sigma}} e_j^{\sigma}(t) + \sigma_j^2 \frac{\partial \pi_w}{\partial w_j} \quad (40)$$

【0 1 1 0】

ただし、ノーターションとしては以下のとおりである。

30

【0 1 1 1】

【数 2 6】

 $\kappa^{\mu}, \kappa^{\sigma}$: それぞれのパラメータのエリジビリティ・トレースの時定数

【0 1 1 2】

(2 - 4 . 具体例)

図4に示した学習システムにおいて、数値シミュレーションを行った結果について以下説明する。

40

【0 1 1 3】

このシミュレーションにおいて、図1に示した2足歩行移動システム1000(2足歩行ロボット)は、脚長が0.2m、両脚の質量がそれぞれ0.5kgとし、胴体が0.1kgであるものとした。さらに、膝関節がないことを考慮して、遊脚を振り出す場合は足先が地面を通過出来るように設定した。

【0 1 1 4】

それぞれの学習パラメータは、以下のとおりである。

【0 1 1 5】

【数 27】

学習率 : $\alpha = 0.04$
 $\beta^H = 0.2$
 $\beta^\sigma = 0.01$
 各エリジビリティの時定数 : $\kappa = 2.0$
 サーボ系のゲイン : $K_p = 5.0$
 $K_d = 0.1$

10

【0116】

また、NGnetの基底関数は、実際にロボットが歩行運動を行う際に必要であると予想される状態空間に格子状に均等に配置することを考え、以下のようにする。

【0117】

【数 28】

$$\begin{aligned} \mathbf{x} &= (x_1, x_2, x_3, x_4)^T \\ &= (\theta_l + \theta_p, \theta_r + \theta_p, \dot{\theta}_l + \dot{\theta}_p, \dot{\theta}_r + \dot{\theta}_p)^T : (12, 6, 12, 6) \end{aligned}$$

【0118】

この結果、計5184 (= 12 × 6 × 12 × 6) 個をそれぞれ、以下の範囲に均等に配置した。

20

【0119】

【数 29】

$$\left(-\frac{\pi}{2} \leq x_1 \leq \frac{\pi}{2}, -3\pi \leq x_2 \leq 3\pi, -\frac{\pi}{2} \leq x_3 \leq \frac{\pi}{2}, -3\pi \leq x_4 \leq 3\pi\right)$$

【0120】

報酬関数は以下の式で表す。

30

【0121】

【数 30】

$$r(\mathbf{x}) = k_H r_H(\mathbf{x}) + k_S r_S(\mathbf{x}) \quad (41)$$

【0122】

ただし、それぞれがロボットの腰の高さに関する項 $r_H(t)$ 、歩行速度に関する項 $r_S(t)$ は、以下の式で表される。

【0123】

【数 31】

40

$$\begin{aligned} r_H(x) &= h_1 - h' - \min(f_l, f_r), \\ r_S &= \max(-1, \min(\dot{x}_1, 1)) \end{aligned}$$

【0124】

ここで、 h_1 はロボットの腰の高さ、 h' は腰の高さのオフセット、 f_l, f_r は左右の脚の高さである。したがって、式(41)の右辺第1項は、ロボットの位置エネルギーに関連する量であり、右辺第2項はロボットの運動エネルギーに関連する量である。

【0125】

以下に説明するシミュレーションでは各パラメータは、 $k_S = 0.06$ 、 $k_H = 0.5$ 、

50

$h' = 0.15$ とした。また、CPGのパラメータは(1-2, 学習システム)で述べたものを用いている。

【0126】

計算機上でのロボット及びCPGのダイナミクスの時間刻みは1 msec、学習システムの時間刻みは10 msecとした。

【0127】

また、シミュレーションにおいて、1学習試行の終了条件は以下のようにした。

【0128】

i) 17700 msec 経過(約100歩の歩行終了後)

ii) 転倒時(ただし、同時に $r = -1$ の報酬を与える)

(2-5, 平地歩行の獲得及び、環境変化に対するロバスト性)

図7は、1試行で獲得した報酬の総和を、試行回数ごとに取った学習曲線を示す図である。図7においては、地面の傾斜0°のときの学習曲線を示している。

【0129】

図7より、学習は約350回で収束しており、定常歩行運動を獲得出来ていることが分かる。

【0130】

図8は、図7の学習曲線に対応する歩行の軌跡を示す図である。

【0131】

図8において、(1)は学習前の歩行軌跡、(2)は600回学習後の歩行軌跡を示す。600回の学習後では、歩幅が大きくなり歩行速度も向上して、良好な歩行軌跡が得られていることがわかる。

【0132】

また、600回学習試行を行うことによって学習した各学習パラメータを用い、数度の傾斜を付けることによって環境を変化させた場合でも、ある程度歩行動作を維持することが可能である。さらに、数回の学習試行を行うことによって、新しい環境に適応することが出来る。これは、行動則の内部状態(ここではCPGの内部状態)と、ロボットの状態が引き込みを行うことによつて、ロバストなリミットサイクルを構成しているからであると考えられる。

【0133】

図9は、図7で獲得した歩行において、CPGの内部状態とロボットの状態間のリミットサイクルを、CPGの内部状態 z_1 と脚角度の時間変化として示す図である。

【0134】

また、図10は、図7で獲得した歩行において、CPGの内部状態とロボットの状態間のリミットサイクルを、脚角度、脚の角速度、CPGの内部状態 z_1 の関係として示す図である。

【0135】

外部からの擾乱に対しても、本発明の制御システムは、周期運動を継続させることが可能なことがわかる。

【0136】

(2-6, 報酬と獲得した運動の関係)

式(41)の報酬関数中の、速度項係数 k_v を変化させた場合の、ロボットの歩行速度の関係を表1に示す。

【0137】

10

20

30

40

【表 1】

表 1: 報酬の速度項係数と歩行速度との関係

報酬の速度項係数 k_s	ロボットの歩行速度 (m/sec)
0.06	0.345
0.1	0.388
0.6	0.419

【0138】

ここで、ロボットの腰の高さに関する項の係数は、前節と同様 $k_H = 0.5$ とした。

【0139】

表 1 より、速度項を増加させるとロボットの歩行速度も増加することが分かり、よってロボットのダイナミクスから構成するようなコントローラを陽に用いることなく、学習の報酬を変化させることによって、ロボットを制御出来ることが確認出来る。

【0140】

(2-7. センサノイズ・時間遅れに対するロバスト性)

図 7 で獲得された歩行を教師信号として学習した各パラメータを初期値として用い、さらに図 4 の学習システムから CPG を取り除いたものを用いて、150 回学習試行を行うことによって、内部状態を持たない行動則によって 2 足歩行運動を獲得した。これと、図 7 の学習によって獲得した歩行運動を用いて、コントローラのセンサノイズ及び時間遅れに対するロバスト性について比較を行った。

【0141】

センサノイズは x_1 、 x_3 に対しては、 $N(0, 0.01)$ 、 x_2 、 x_4 に対しては、 $N(0, 0.09)$ を用い、時間遅れは 20 msec としてシミュレーションを行った。

【0142】

図 11 は、センサノイズ・時間遅れに対するシミュレーション結果を示す図である。図 11 において、(1) は CPG 有り、(2) は CPG 無しコントローラで構成された歩行を示す。また、図 11 において、(a) は通常の条件での歩行、(b) はセンサノイズのある状態での歩行、(c) は時間遅れがある場合の歩行であり、また、図 11 中で、“ ” はロボットの進行方向を表している。

【0143】

CPG を持たない行動則で構成された歩行は、ノイズ及び時間遅れのどちらの場合についても歩行動作を保つ事は出来なかったが、図 4 に示した学習システムでは、ノイズ及び時間遅れがある場合でも歩行が可能であることがわかる。

【0144】

よって内部状態を持つ行動則を構成することによって、センサノイズや時間遅れに対してロバストなコントローラを構成出来ることが分かる。

【0145】

(3. 正規化ガウス関数ネットワークによる関数近似)

2-3-1 で述べた価値関数、フィードバック制御器を表現するために用いた、正規化ガウス関数ネットワークについて、以下説明する。

【0146】

NGnet は 3 層のネットワークで構成されており、中間素子は正規化ガウス関数である。

入力ベクトル $x = (x_1, \dots, x_n)^T$ に対して、 k 番目のユニットの活性化関数は、以下の式のようになる。

【0147】

10

20

30

40

【数 3 2】

$$\phi_k(\mathbf{x}) = e^{-\frac{1}{2}\|M_k(\mathbf{x}-\mathbf{c}_k)\|^2} \quad (42)$$

【0 1 4 8】

ただし、 c_k は活性化関数の中心であり、 M_k は活性化関数の形状を決定する行列である。ここで、活性化関数 $\phi_k(x)$ を各点で総和が1になるように以下の式(43)のように正規化したものを、基底関数 $b_k(x)$ とする。

【0 1 4 9】

【数 3 3】

$$b_k(\mathbf{x}) = \frac{\phi_k(\mathbf{x})}{\sum_{l=1}^K \phi_l(\mathbf{x})} \quad (43)$$

10

【0 1 5 0】

ただし、 K は基底関数の個数である。

【0 1 5 1】

このような正規化を行うことによって、中心点 c_k が密に配置されている部分では、 $b_k(x)$ は局所的な基底関数となり、 c_k の分布の端の部分では $b_k(x)$ はシグモイド関数のような大域的な基底関数になる。

20

【0 1 5 2】

ネットワークの出力は、基底関数と重みの内積によって以下の式(44)ようになる。

【0 1 5 3】

【数 3 4】

$$y_k(\mathbf{x}) = \sum_{k=1}^K w_k b_k(\mathbf{x}) \quad (44)$$

【0 1 5 4】

この出力が、図4の正規化処理部1030の出力となる。

30

【0 1 5 5】

[実施の形態1の変形例]

以上の説明では、図3(3)の構成による制御について説明した。以下では、実施の形態1の変形例として、図3(2)の構成による制御について説明する。

【0 1 5 6】

図12は、図3(2)の構成に相当するシステムであって、制御器と制御対象を含めたシステム全体の構成を示す図である。

【0 1 5 7】

図12において、状態観測器2002は、状態観測器のダイナミクス2004と、方策勾配法(強化学習)に基づいた強化学習器2006によって構成される。

40

【0 1 5 8】

状態観測器2002中の強化学習器2006は、以下に説明するとおり、制御対象の観測出力 y と、状態観測器のダイナミクス2004に基づく出力と、制御器2010の制御出力 u とに基づいて、学習器出力 U を出力する。出力関数処理部2030は、状態観測器2002からの推定状態に基づいて、状態観測器2002の出力を報酬演算部2020に与える。報酬演算部2020は、状態観測器2002の出力と観測対象からの観測出力 y と学習器出力 U とに基づいて、報酬を計算し、強化学習器2006に与える。

【0 1 5 9】

(方策勾配法を用いた状態観測器の学習)

50

以下では、実施の形態 1 の変形例の状態観測器 2 0 0 2 の構造について説明する。

【 0 1 6 0 】

状態の推定値を、 x の頭部に “ ^ ” を付加して表現 (= x_i) (以下、本文中では 「 x ハット 」 と呼ぶ) したとき、つぎのような状態観測器を考える。

【 0 1 6 1 】

【 数 3 5 】

$$\begin{aligned} \dot{\hat{x}} &= \mathbf{f}(\hat{x}, \mathbf{u}(\hat{x})) + \mathbf{U}(\hat{x}, \mathbf{u}(\hat{x}), \mathbf{y}) \\ &= \mathbf{g}(\hat{x}, \mathbf{u}(\hat{x}), \mathbf{U}) \end{aligned}$$

10

【 0 1 6 2 】

ここではまず、通常のオブザーバやカルマンフィルタ同様、対象のダイナミクス $f(x, u)$ は既知または学習によって獲得可能であるとし、対象システムの観測出力 y を基にして、現在の推定状態 x ハットと制御出力 u から、推定状態を真の状態にどのように近づけるべきかを方策勾配法を用いて学習する。

【 0 1 6 3 】

ここでは学習器の目的を、状態観測器の出力 y ハット (y の頭部に “ ^ ” を付加したもの) と対象システムの出力 y との誤差を最小にすることとする。

20

【 0 1 6 4 】

よって、報酬演算部 2 0 2 0 により演算される報酬関数は次のようになる。

【 0 1 6 5 】

【 数 3 6 】

$$r(t) = -(\mathbf{y}(t) - \hat{\mathbf{y}}(t))^T \mathbf{Q}(\mathbf{y}(t) - \hat{\mathbf{y}}(t)) - \mathbf{U}^T(t) \mathbf{R} \mathbf{U}(t)$$

【 0 1 6 6 】

ただし、 \mathbf{Q} , \mathbf{R} は報酬関数の形を決めるパラメータである。この結果、学習器は状態観測器のダイナミクス 2 0 0 4 への以下のようなノーターションのフィードバック入力 \mathbf{U} を獲得することになる。

30

【 0 1 6 7 】

【 数 3 7 】

$$\mathbf{U}(\hat{x}, \mathbf{u}, \mathbf{y})$$

【 0 1 6 8 】

ここで、フィードバック入力 \mathbf{U} は次の確率分布により表現される。

【 0 1 6 9 】

【 数 3 8 】

$$\pi(U_j, \hat{x}, \mathbf{y}; \mathbf{w}) = \frac{1}{\sigma_j \sqrt{2\pi}} \exp\left(\frac{-(U_j - \mu_j)^2}{2\sigma_j^2}\right)$$

40

【 0 1 7 0 】

したがって、 j 番目の出力の実現値 U_j は、以下の式により与えられる。

【 0 1 7 1 】

【数 3 9】

$$U_j(t) = \mu_j(t) + \sigma_j n_j(t)$$

【0172】

ただし、 $n_j(t) \sim N(0, 1)$ であり、 $N(0, 1)$ は、上述のとおり、平均 0、分散 1 の正規分布を表す。フィードバック入力 U を生成する確率分布の更新は、二足歩行運動の学習の場合と同様に行われる。

【0173】

このような構成によっても、周期運動に対し状態観測器の学習を容易とすることが可能な動的制御装置およびこのような動的制御装置を用いた 2 足歩行移動装置を提供することができる。

10

【0174】

[実施の形態 2 (多自由時系への適用)]

(5 . 1) 5 リンク 2 足歩行ロボットモデル

実施の形態 2 では、図 4 で構成した学習システムを、多自由度系へ適用した構成について説明する。

【0175】

実施の形態 1 で述べた 3 リンク 2 足歩行ロボットモデルに対する学習結果から、学習手法に方策勾配法を用いることにより、物理系の状態のみを考慮することで、望ましい行動則を獲得出来ることが分かる。

20

【0176】

図 1 3 は、実施の形態 2 で扱う 5 リンク 2 足歩行ロボットおよびそのシミュレーションモデルを説明するための図である。

【0177】

実施の形態 1 の結果を基にして、以下では、図 1 3 (a) に示した膝関節を持つ 5 リンクの実ロボットのシミュレーションモデル (図 1 3 (b)) に対して、実施の形態 1 で提案した手法を適用し、歩行運動の獲得を行なう。

【0178】

図 1 3 (b) に示すとおり、実施の形態 2 のシミュレーションモデルでは、右下腿 1 0 r d (リンク 1)、右上腿 1 0 r u (リンク 2)、腰部 4 0 (リンク 2)、左上腿 1 0 l u (リンク 4) および左下腿 1 0 l d (リンク 5) の 5 リンク系である。

30

【0179】

図 1 3 (b) に示す 2 足歩行ロボットも、動的制御装置 1 0 0 と、動的制御装置 1 0 0 の上部に設けられる胴部 4 0 と、動的制御装置 1 0 0 により駆動制御される脚部とを備える。脚部は、右上腿 1 0 r u と、右下腿 1 0 r d と、左上腿 1 0 l u と、左下腿 1 0 l d とを有する。

【0180】

図 1 3 (b) に示す 2 足歩行ロボットは、各脚において、右足および左足のそれぞれ接地面近傍に設けられるセンサ 2 0 r および 2 0 l と、右膝部分に設けられるセンサ 2 0 r k と、左膝部分に設けられるセンサ 2 0 l k とを備えるものとする。一方、胴部 4 0 には、センサ (図示せず) が設けられる。動的制御装置 1 0 0 は、実施の形態 1 と同様に、駆動部 1 0 8 により、右上腿 1 0 r u および左上腿 1 0 l u を駆動する。また、左右の膝にも、図示しない駆動部が設けられ、動的制御装置 1 0 0 からのトルク信号に基づいて、右下腿 1 0 r d および左下腿 1 0 l d を駆動するものとする。

40

【0181】

センサ 2 0 r k は、胴部 4 0 の中心線 4 に対する右上腿 1 0 r u の角度 θ_{hip} 、角速度 $d \theta_{hip} / dt$ という情報を検出し、また、センサ 2 0 l k は、中心線 4 に対する左上腿 1 0 l u の角度 θ_{hip} 、角速度 $d \theta_{hip} / dt$ という情報とを検出し、それぞれ、動的制御装置 1 0 0 に通知する。さらに、センサ 2 0 r は、右上腿 1 0 r u の延長線 6 r に対す

50

る右脚の角度 θ_{knee} 、角速度 $d\theta_{knee}/dt$ という情報と右足の接地状態とを検出し、また、センサ 201 は、左上腿 101u の延長線 61 に対する左脚の角度 θ'_{knee} 、角速度 $d\theta'_{knee}/dt$ という情報と左足の接地状態とを検出し、それぞれ、動的制御装置 100 に通知する。

【0182】

また、胸部 40 に設けられるセンサは、鉛直方向 2 に対する胸部 40 のピッチ角 θ_p 、角速度 $d\theta_p/dt$ という情報とを検出し、それぞれ、動的制御装置 100 に通知する。

【0183】

5リンク2足歩行ロボットモデルの物理パラメータを表1に示す。

【0184】

【表2】

	リンク1	リンク2	リンク3	リンク4	リンク5
質量 [kg]	0.15	0.64	2.0	0.64	0.15
長さ [m]	0.20	0.20	0.01	0.20	0.20
慣性モーメント ($\times 10^{-4}$) [kg·m ²]	1.40	6.89	1.0	6.89	1.40
各関節から重心まで の長さ [m]	0.1	0.12	0.0	0.12	0.1

【0185】

(5.2) 多自由度系に対する学習システム

図13(b)に示す5リンク2足歩行ロボットモデルに対する学習システムについて、以下さらに詳しく説明する。

【0186】

動的制御装置 100 では、膝関節の状態は学習に用いず、実施の形態1の図4と同様に、腰関節の状態 (θ_{hip} 、 θ'_{hip} 、 $d\theta_{hip}/dt$ 、 $d\theta'_{hip}/dt$) 及びピッチ角の状態 (θ_p 、 $d\theta_p/dt$) のみに関する状態変数を用いて学習を行う。

【0187】

一方、膝関節については、接地情報と腰関節の状態を基に膝関節の目標関節角を切替える状態マシン 1040 を制御器として用いる。この状態マシン 1040 については後述する。

【0188】

図14は、動的制御装置 100 における実施の形態2の学習システムの構成を示す図である。なお、図4に示した実施の形態1の動的制御装置 100 と同一部分は、同一符号で示す。ただし、図14においては、フィードバック制御器 1022' は、フィードバック制御器 1022 の機能と出力飽和处理部 1024 の機能とを併せて有するものとし、価値関数処理部 1032' は、価値関数処理部 1032 の機能と正規化处理部 1030 の機能も併せて有するものとする。

【0189】

左上腿 101u および右上腿 10ru へのトルクは、PDサーボ処理部 1028 から与えられる。

【0190】

したがって、価値関数処理部 1032'、フィードバック制御器 1022'、CPG 処理部 1026 および PDサーボ処理部 1028 の動作は、駆動する対象が異なるのみで、その基本的な動作は、実施の形態1と同様である。

【0191】

一方、膝関節へのトルク入力は、状態マシン 1040 によって決定される目標関節角を用いた PDサーボ処理部 1042 により与えられる。

10

20

30

40

50

【 0 1 9 2 】

(状態マシン 1040 による膝関節制御器)

図 15 は、状態マシン 1040 の動作を説明するための概念であり、図 15 (a) は、右膝の状態を、図 15 (b) は左膝の状態をそれぞれ示す。

【 0 1 9 3 】

状態は、膝屈曲 膝伸長 膝屈曲 膝伸長との状態遷移を繰り返す。

【 0 1 9 4 】

たとえば、右膝では、前半の膝屈曲 膝伸長は、脚の「振り状態」であり、目標角度はそれぞれ θ_1 および θ_2 である。右膝の後半の膝屈曲 膝伸長は、脚の「立ち状態」であり、目標角度はそれぞれ θ_3 および θ_4 である。

10

【 0 1 9 5 】

一方で、左膝では、前半の膝屈曲 膝伸長は、脚の「立ち状態」であり、目標角度はそれぞれ θ_3 および θ_4 である。左膝の後半の膝屈曲 膝伸長は、脚の「振り状態」であり、目標角度はそれぞれ θ_1 および θ_2 である。

【 0 1 9 6 】

腰関節の角度 θ_{hip} 、 $\dot{\theta}_{hip}$ (以下、総称するときは、「 θ_{hip} 」) および接地情報を用いて、図 15 に示す状態マシンにより、膝関節の目標角度 (θ_{knee} ハット、 $\dot{\theta}_{knee}$ ハット: 総称するときは、「 θ_{knee} ハット」) を決定し、以下に示す PD サーボにより膝関節へのトルク u_{knee} を出力する。

【 0 1 9 7 】

【数 4 0】

20

$$u_{knee} = K_p (\hat{\theta}_{knee} - \theta_{knee}) - K_d \dot{\theta}_{knee}$$

【 0 1 9 8 】

ここで、 K_p は位置ゲイン、 K_d は速度ゲインである。ただし、シミュレーションでは、 $K_p = 12.0$ 、 $K_d = 0.15$ とした。

【 0 1 9 9 】

図 15 に示すように、状態マシンでは 4 つの目標角 (θ_1 、 θ_2 、 θ_3 、 θ_4) を設定した条件に応じて切り換える。ここでは、 $\theta_1 = 1.11$ 、 $\theta_2 = 0.56$ 、 $\theta_3 = 0.52$ 、 $\theta_4 = 0.26$ とした。膝を曲げた状態から伸ばした状態への遷移は腰関節を用いた以下の条件に基づいて行なわれる。

30

【 0 2 0 0 】

【数 4 1】

$$\theta_{hip}^l - \theta_{hip}^r < b \quad \text{または} \quad \theta_{hip}^r - \theta_{hip}^l < b$$

【 0 2 0 1 】

膝を伸ばした状態から曲げた状態への遷移は足裏の接地情報により行われる。シミュレーションでは、たとえば、 $b = 0.15$ とする。

40

【 0 2 0 2 】

(床反力モデル)

以下では、さらに、シミュレーションで用いる床反力モデルについて説明する。

【 0 2 0 3 】

x 、 y は脚端の位置を表し、 x_g 、 y_g は接地点とすると、接地時の床反力は以下の式でモデル化される。

【 0 2 0 4 】

【数 4 2】

$$F_x = k_x(x_g - x) - b_x \dot{x}$$

$$F_y = k_y(y_g - y) - b_y \dot{y}$$

if $y \geq y_g$

$$F_x = 0, F_y = 0$$

【0 2 0 5】

ここで、 F_x 、 F_y は水平方向、垂直方向の床反力である。

10

【0 2 0 6】

以下で説明するシミュレーションでは、上記の式のそれぞれの係数は $k_x = 3000$ 、 $b_x = 10$ 、 $k_y = 30000$ 、 $b_y = 100$ とする。また、床反力が $F_x > \mu F_y$ を満たすときに、足裏が床面を滑ると定義し、ここでの μ は静摩擦係数であり、 $\mu = 1.0$ としている。

【0 2 0 7】

(5.3)シミュレーション

以下、実施の形態2でのシミュレーション結果を説明する。

【0 2 0 8】

本シミュレーションで用いるCPGのパラメータは、ロボットモデルの脚長の変化を考慮して、図4の3リンク2足歩行ロボットモデルで用いた特性と比較して、低い周波数が必要となるので、各パラメータを以下のようにする。

20

【0 2 0 9】

$$= 0.4, \quad \dot{\quad} = 0.8, \quad = 1.3, \quad = 2.0, \quad z_0 = 0.1$$

正規化ガウス関数ネットワークNGnetの基底関数については、それぞれの状態変数に対して格子状に、 $x = (x_1, x_2, x_3, x_4)^T$ を以下のように配置する。

【0 2 1 0】

【数 4 3】

$$\begin{aligned} x &= (x_1, x_2, x_3, x_4)^T \\ &= (\theta_{hip}^l + \theta_p, \theta_{hip}^r + \theta_p, \dot{\theta}_{hip}^l + \dot{\theta}_p, \dot{\theta}_{hip}^r + \dot{\theta}_p)^T : (16, 10, 16, 10) \end{aligned}$$

30

【0 2 1 1】

基底関数は、25600個をそれぞれの以下の範囲に均等に配置した。

【0 2 1 2】

【数 4 4】

$$\left(-\frac{\pi}{3} \leq x_1 \leq \frac{\pi}{3}, -2.5\pi \leq x_2 \leq 2.5\pi, -\frac{\pi}{3} \leq x_3 \leq \frac{\pi}{3}, -2.5\pi \leq x_4 \leq 2.5\pi\right)$$

40

【0 2 1 3】

報酬関数に関しては、式(41)を使用するが、物理パラメータの変更にともない、 $k_s = 0.06$ 、 $k_H = 0.5$ 、 $h' = 0.3$ とした。

【0 2 1 4】

その他の学習パラメータ等については、式(38)の (は上付き) = 0.02以外は、実施の形態1と同様とする。

【0 2 1 5】

また、1学習試行の終了条件についても、実施の形態1と同様以下の条件を用いる。

50

【0216】

- ・ 17700 msec 経過
 - ・ 転倒時(ただし、同時に $r = -1$ の報酬を与える)
- (5.4) 学習結果

実施の形態2のシミュレーションにおいても、3リンク2足歩行ロボットでの学習と同様、ロボットが10試行連続で、転倒せずに歩き続けたとき、2足歩行運動を獲得したと定義する。10回シミュレーションを行った結果全ての試行で2足歩行運動の獲得に成功した。また、運動獲得に必要な平均試行回数は899回であった。

【0217】

図16は、学習過程の一例の試行回数と獲得報酬の総和の関係を示す図であり、図17は、学習前、学習後の歩行軌跡を示す図である。

10

【0218】

図17において、(1)は学習前の歩行軌跡であり、(2)は1500回学習後の歩行軌跡を示している。

【0219】

この結果は、2足歩行ロボットの自由度が増加し、状態空間が高次元となった場合でも、少ない状態変数だけを観測することで、少ない試行回数で目的の周期運動を獲得できることがわかる。

【0220】

図16の結果からすると、5リンクの2足歩行ロボットに対して、強化学習を用いた従来例(たとえば、文献:Y.Nakayama, M.Sato, and S. Ishii. Reinforcement Learning for biped robot. In Proceedings of the 2nd International Symposium on Adaptive Motion of Animals and Machines, pp. Thp-11-5, 2003.)と比較しても、かなり少ない試行回数で歩行運動を獲得できている。

20

【0221】

以上のような構成によって、2足歩行ロボットの自由度が増加し、状態空間が高次元となった場合でも、周期運動に対し状態観測器の学習を容易とすることが可能な動的制御装置およびこのような動的制御装置を用いた2足歩行移動装置を提供することができる。

【0222】

今回開示された実施の形態はすべての点で例示であって制限的なものではないと考えられるべきである。本発明の範囲は上記した説明ではなくて特許請求の範囲によって示され、特許請求の範囲と均等の意味および範囲内でのすべての変更が含まれることが意図される。

30

【図面の簡単な説明】

【0223】

【図1】本発明の動的制御装置を用いた2足歩行移動システム1000の一例を示す概念図である。

【図2】図1に示した動的制御装置100の構成を示すブロック図である。

【図3】動的行動則を説明するための概念図である。

【図4】演算処理部102の行う処理を示す機能ブロック図である。

40

【図5】神経振動子モデルによるCPGを示す概念図である。

【図6】CPGの出力qを構成する変数 z_1 、 z_2 の波形を示す図である。

【図7】1試行で獲得した報酬の総和を、試行回数ごとに取った学習曲線を示す図である。

【図8】図7の学習曲線に対応する歩行の軌跡を示す図である。

【図9】CPGの内部状態とロボットの状態の間のリミットサイクルを、CPGの内部状態 z_1 と脚角度の時間変化として示す図である。

【図10】CPGの内部状態とロボットの状態の間のリミットサイクルを、脚角度、脚の角速度、CPGの内部状態 z_1 の関係として示す図である。

【図11】センサノイズ・時間遅れに対するシミュレーション結果を示す図である。

50

【図12】制御器と制御対象を含めたシステム全体の構成を示す図である。

【図13】実施の形態2で扱う5リンク2足歩行ロボットおよびそのシミュレーションモデルを説明するための図である。

【図14】動的制御装置100における実施の形態2の学習システムの構成を示す図である。

【図15】状態マシン1040の動作を説明するための概念である。

【図16】学習過程の一例の試行回数と獲得報酬の総和の関係を示す図である。

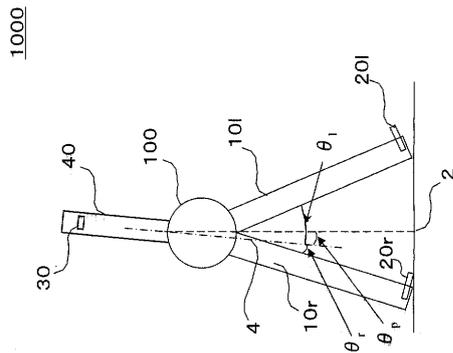
【図17】学習前、学習後の歩行軌跡を示す図である。

【符号の説明】

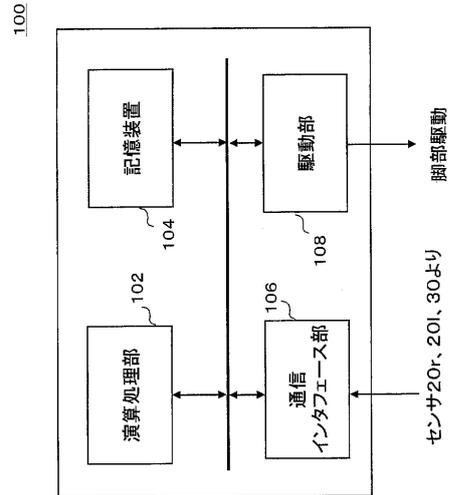
【0224】

10r, 10l 脚部、20r, 20l, 30 センサ、40 胴部、100 動的制御装置、102 演算処理部、104 記憶装置、106 通信インタフェース、108 駆動部、1000 2足歩行移動システム。

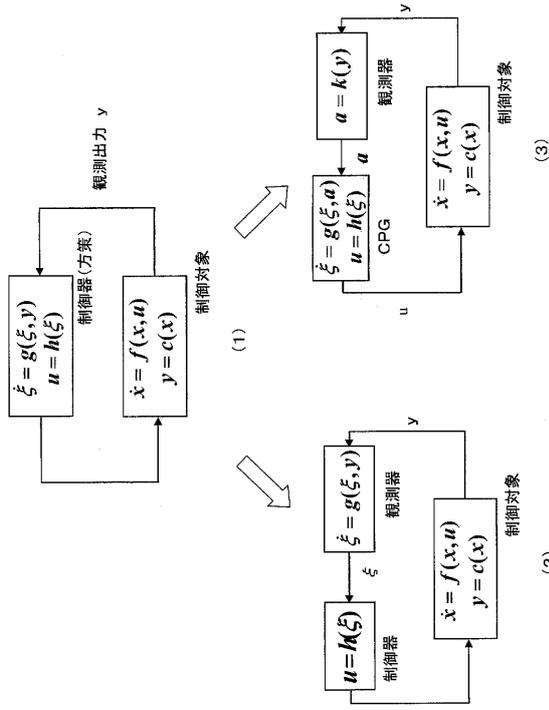
【図1】



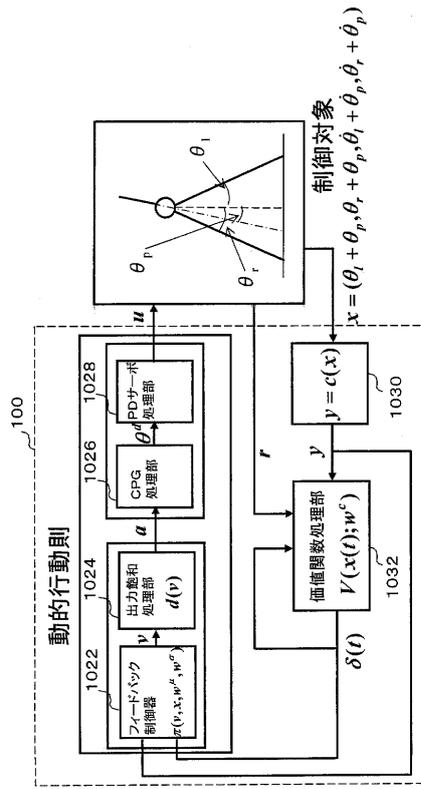
【図2】



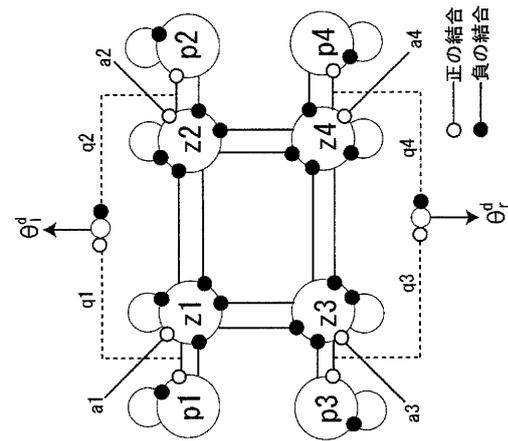
【 図 3 】



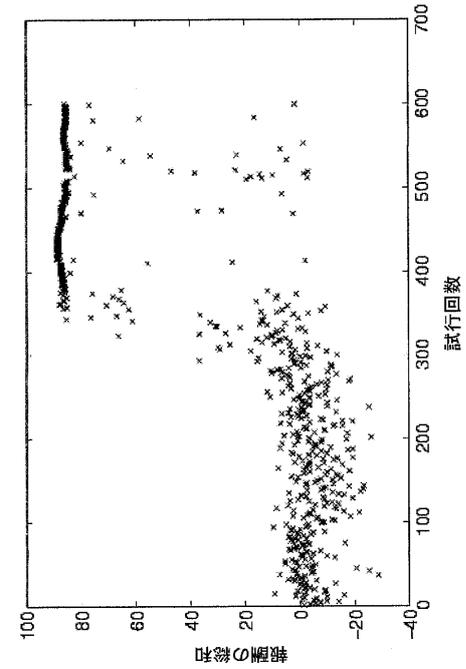
【 図 4 】



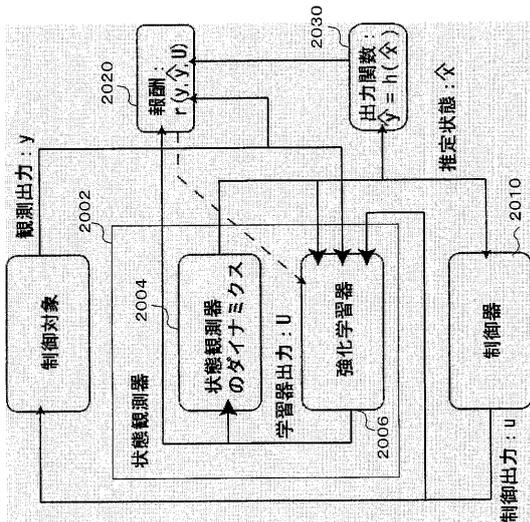
【 図 5 】



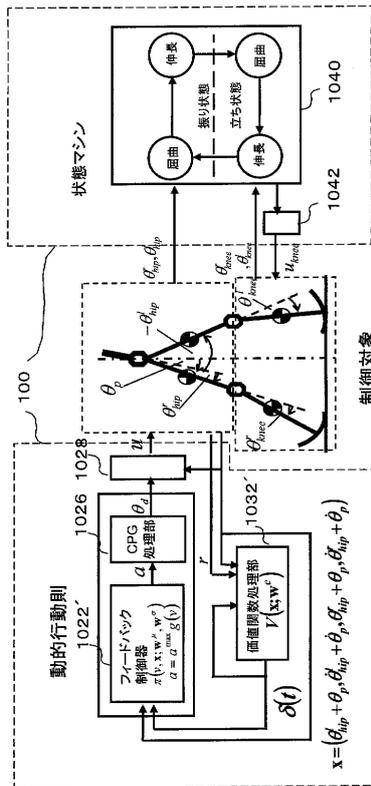
【 図 7 】



【図 12】

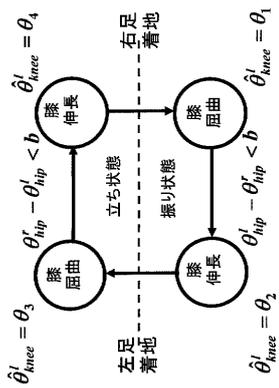


【図 14】



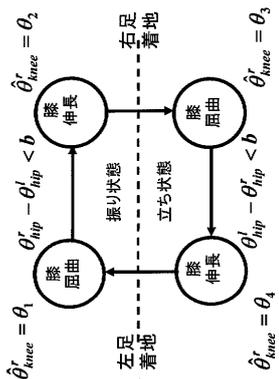
【図 15】

左膝



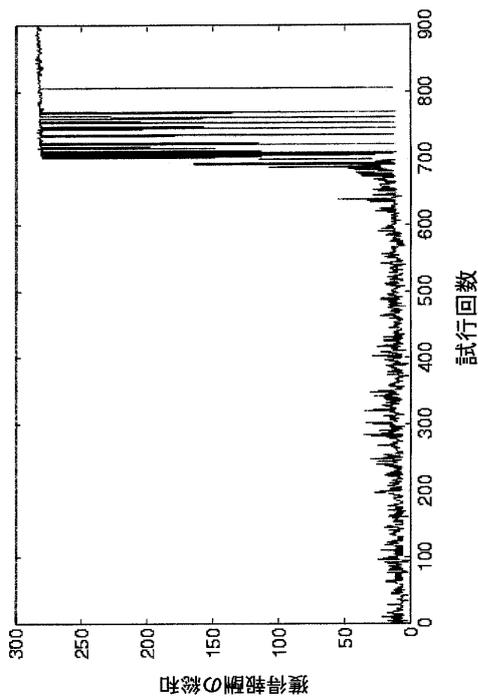
(b)

右膝

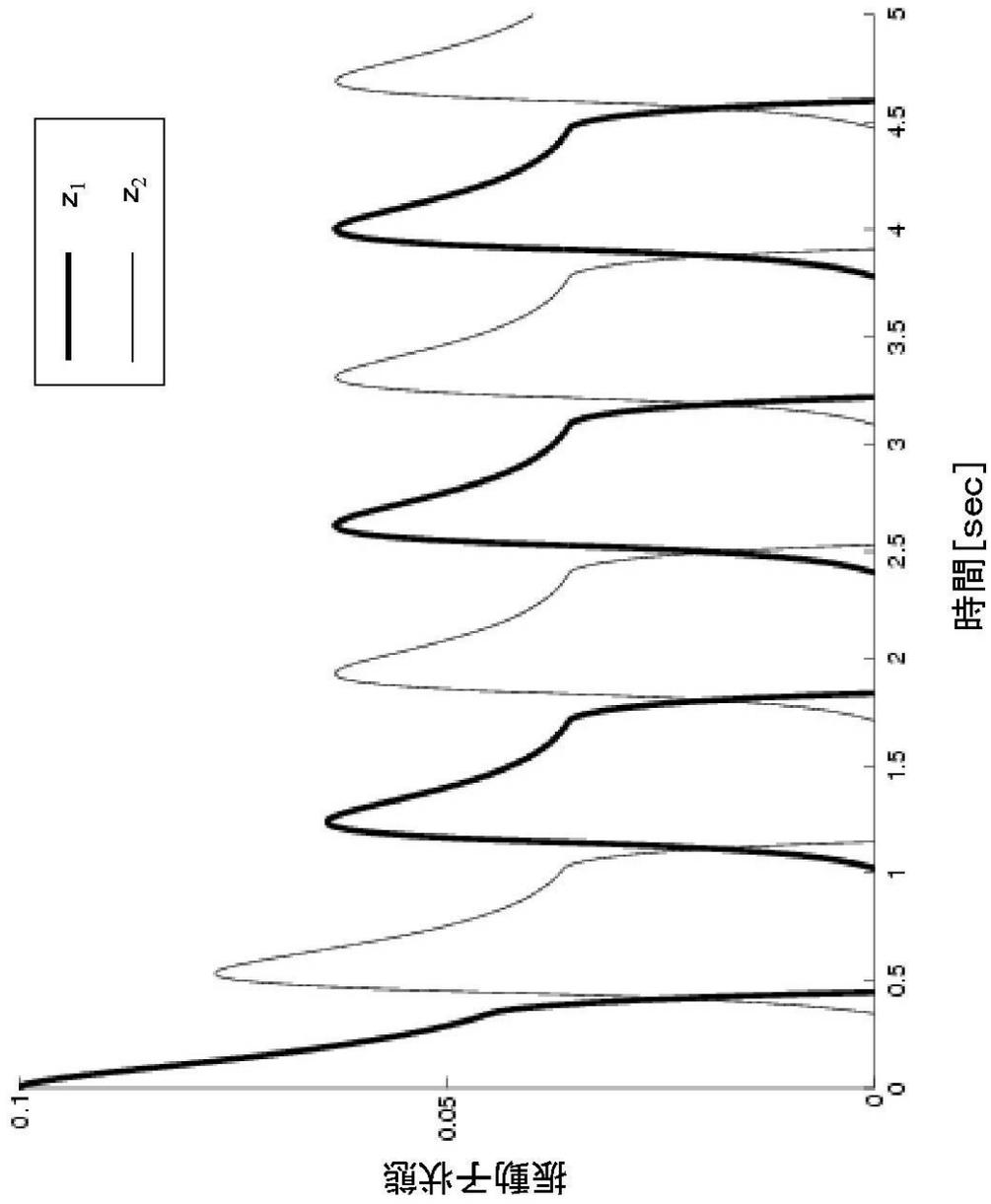


(a)

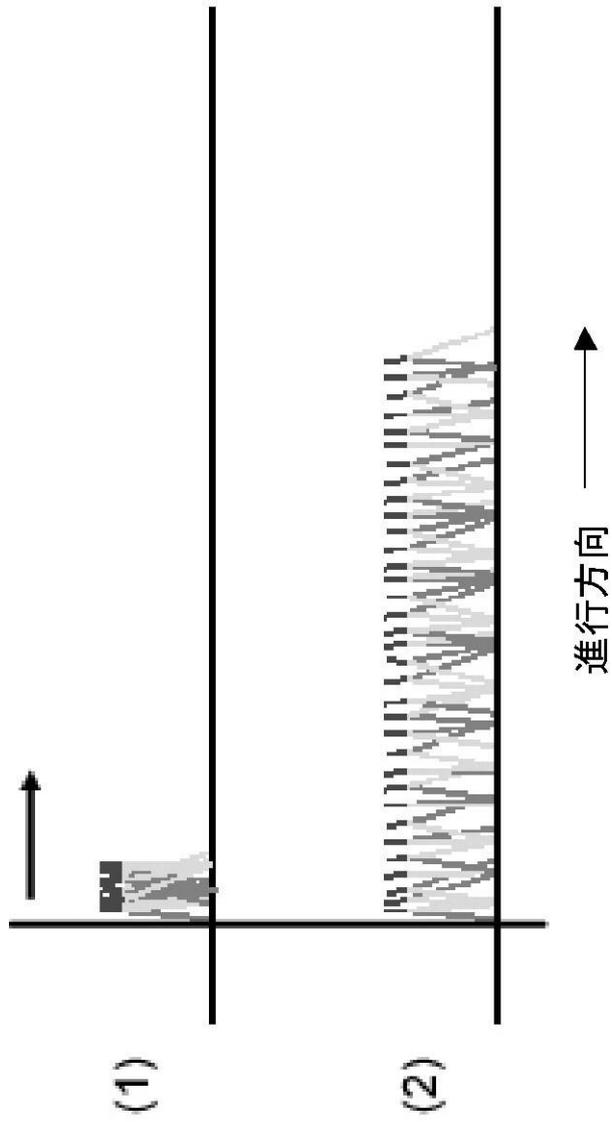
【図 16】



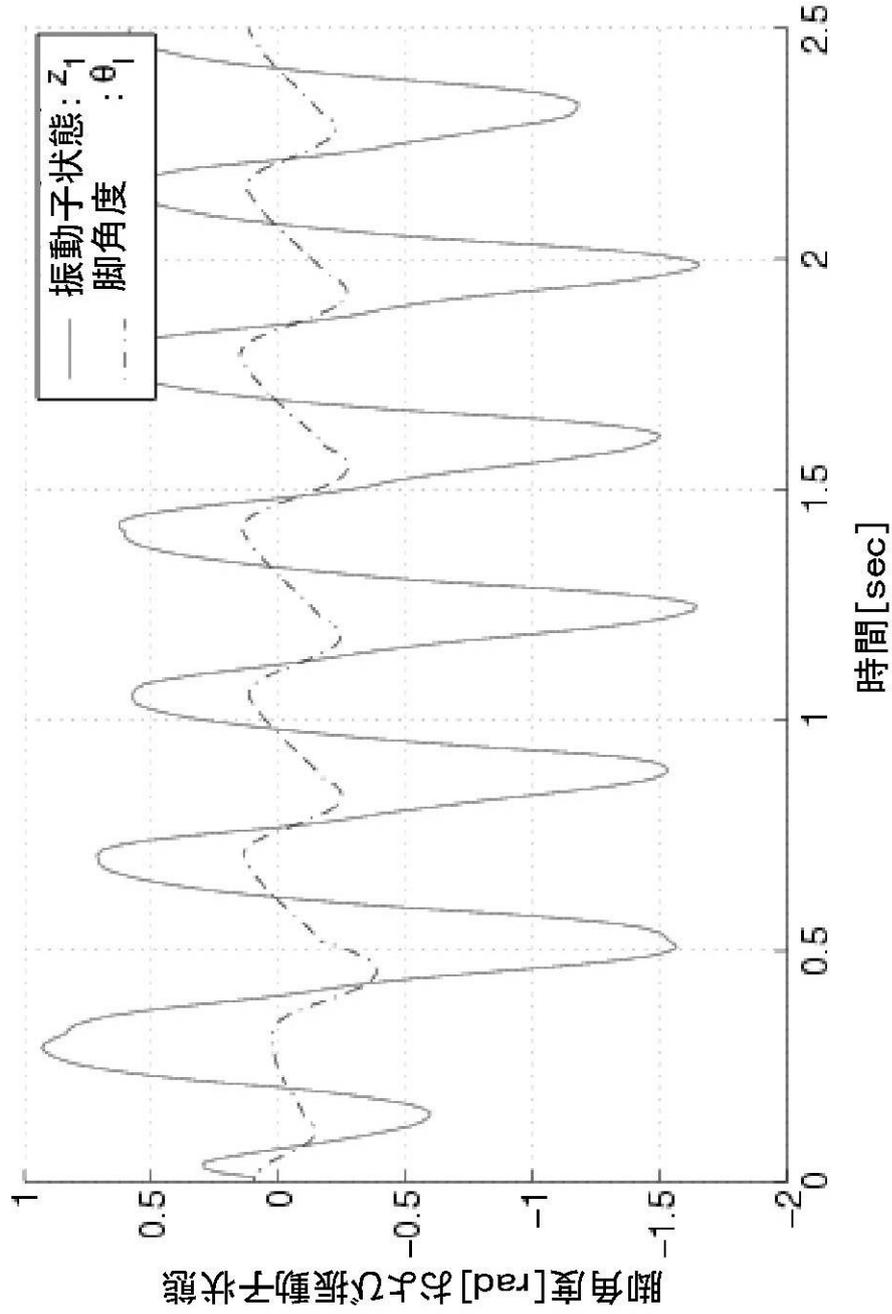
【図6】



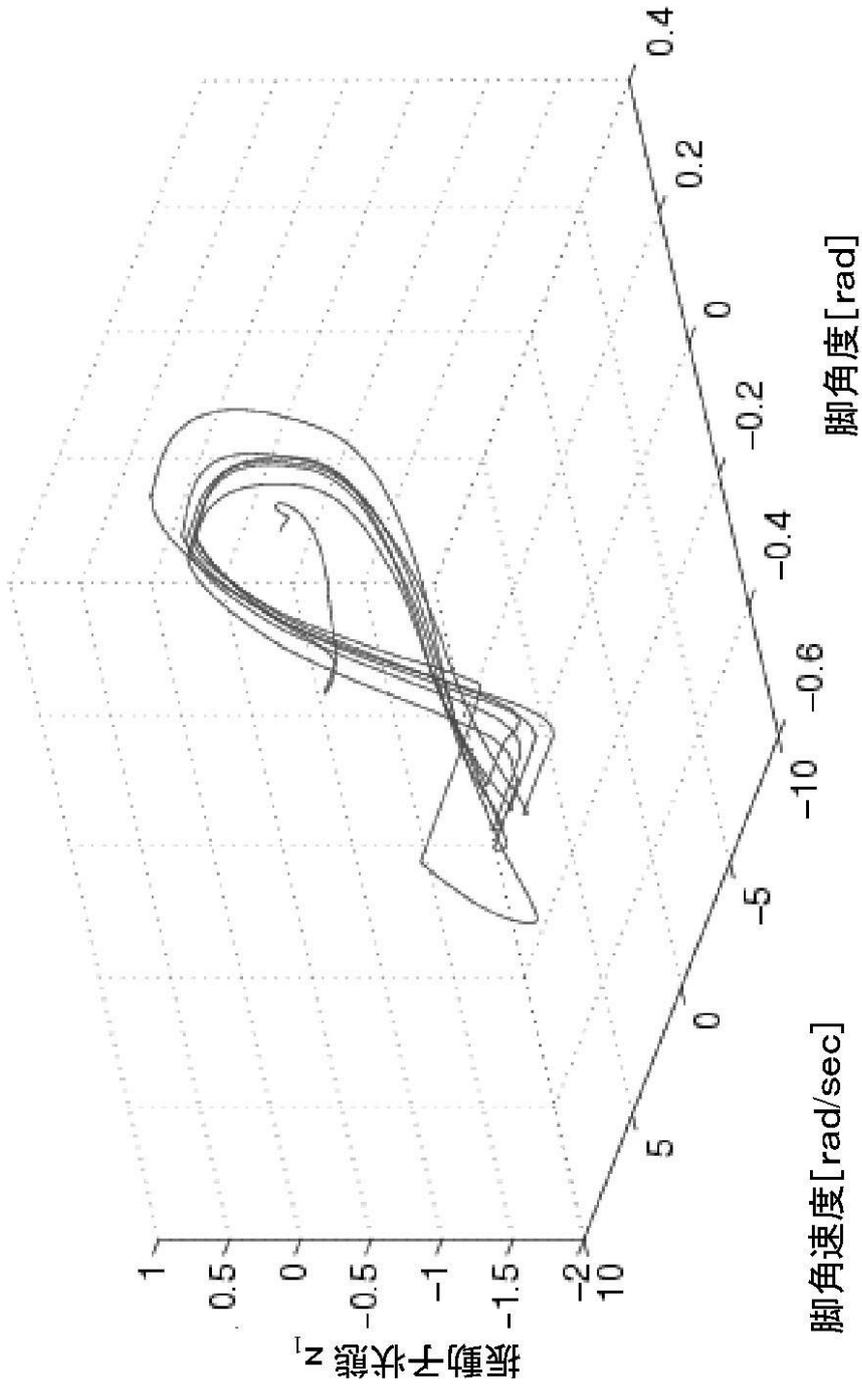
【 図 8 】



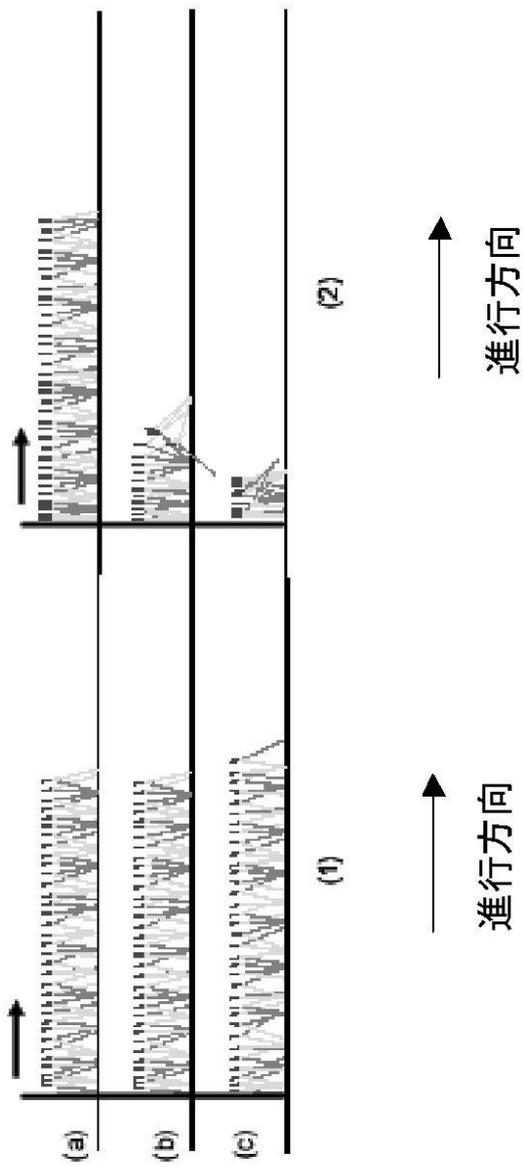
【図9】



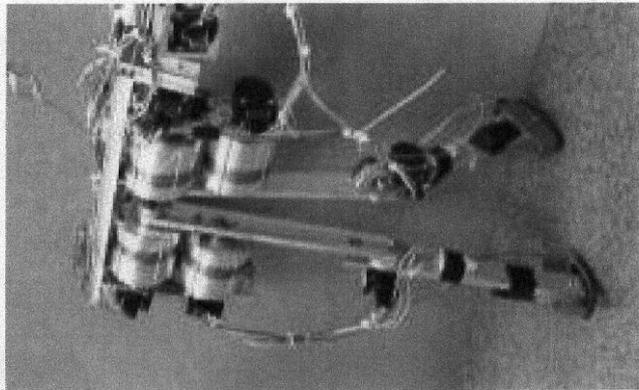
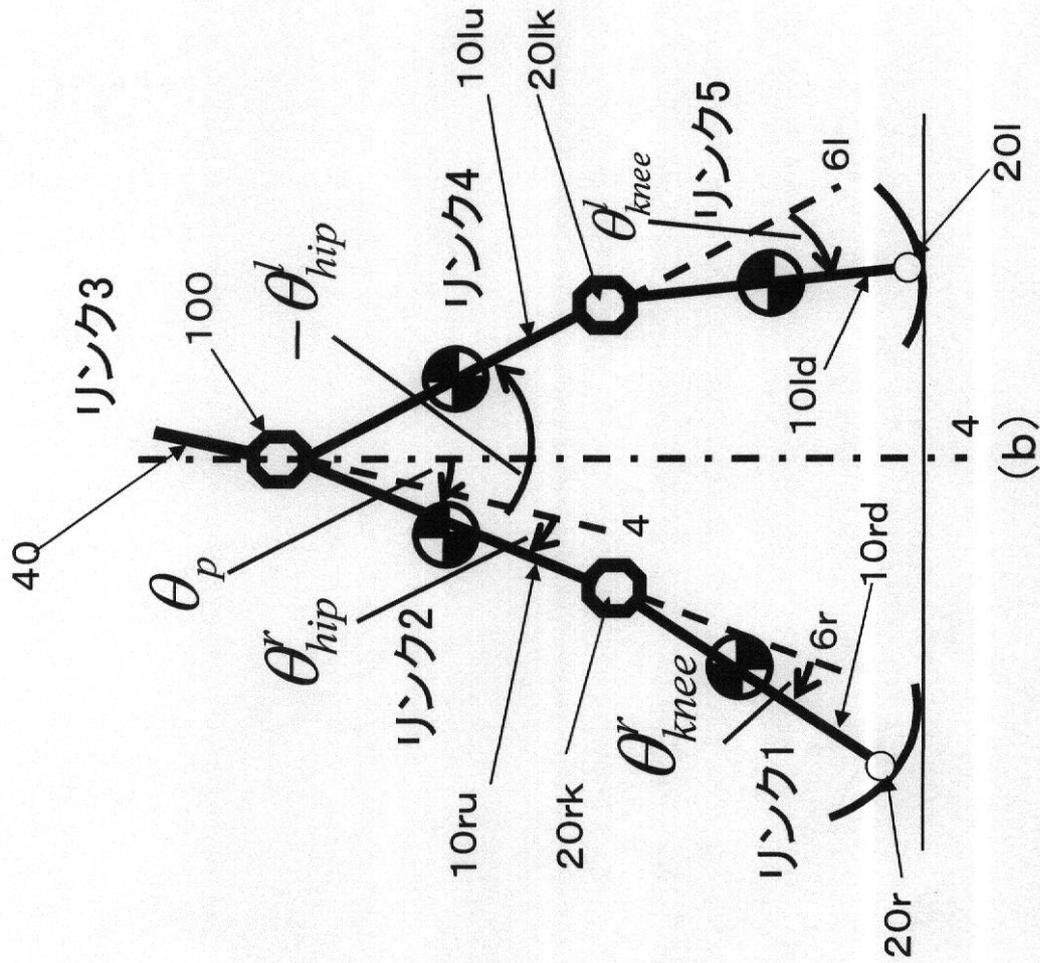
【図10】



【図 11】

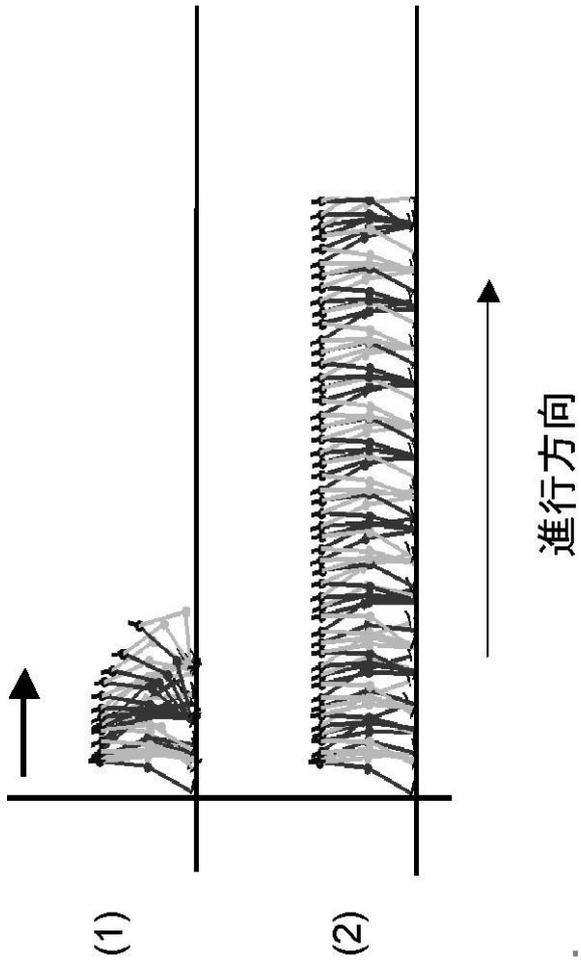


【 図 13 】



(a)

【 図 17 】



フロントページの続き

- (74)代理人 100109162
弁理士 酒井 將行
- (72)発明者 森本 淳
京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
- (72)発明者 松原 崇充
京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
- (72)発明者 佐藤 雅昭
京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
- (72)発明者 中西 淳
京都府相楽郡精華町光台二丁目2番地2 株式会社国際電気通信基礎技術研究所内
- (72)発明者 遠藤 玄
東京都品川区北品川6丁目7番35号 ソニー株式会社内

審査官 植村 森平

- (56)参考文献 森健、中村泰、石井信、二足歩行運動に対する方策勾配法に基づいた強化学習法、電子情報通信学会技術研究報告、日本、社団法人 電子情報通信学会、2004年 3月19日、Vol. 103, No. 734, pp. 73-78
中村泰、佐藤雅昭、石井信、神経振動子ネットワークを用いたリズム運動に対する強化学習法、電子情報通信学会論文誌D-2、日本、社団法人 電子情報通信学会、2004年 3月 1日、Vol. J87-D-2, No. 3, pp.893-902
松原崇充、森本淳、中西淳、佐藤雅昭、銅谷賢治、方策勾配法を用いた動的行動則の獲得：2足歩行運動への適用、電子情報通信学会技術研究報告、日本、社団法人 電子情報通信学会、2004年 1月27日、Vol. 103, No. 602, pp. 53-58
Marc H. Raibert, LEGGED ROBOTS, Communications of the ACM, 米国, Association for Computing Machinery, 1986年 6月, Vol. 29, issue 6, pp. 499-514
中村泰、石井信、佐藤雅昭、神経振動子ネットワークを用いた強化学習法による歩行運動の獲得、電気情報通信学会技術研究報告、日本、社団法人 電子情報通信学会技術研究報告、2002年 3月18日、Vol. 101, No. 735, pp. 183-190
森健、吉本潤一郎、石井信、確率の方策勾配法に基づくactor-critic法と連続システムの制御への応用、電子情報通信学会技術研究報告、日本、社団法人 電子情報通信学会、2003年 3月19日、Vol. 102, No. 731, pp. 137-142

(58)調査した分野(Int.Cl., DB名)

B25J 5/00
JSTPlus(JDreamII)
Cinii