

(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2013-252247
(P2013-252247A)

(43) 公開日 平成25年12月19日(2013.12.19)

(51) Int.Cl.
A63F 13/10 (2006.01)

F I
A63F 13/10

テーマコード(参考)
2C001

審査請求 未請求 請求項の数 8 O L (全 19 頁)

(21) 出願番号 特願2012-128918 (P2012-128918)
(22) 出願日 平成24年6月6日(2012.6.6)

(71) 出願人 504238806
国立大学法人北見工業大学
北海道北見市公園町165番地
(74) 代理人 100081271
弁理士 吉田 芳春
(74) 代理人 100159628
弁理士 吉田 雅比呂
(74) 代理人 100162189
弁理士 堀越 真弓
(72) 発明者 前田 康成
北海道北見市公園町165番地 国立大学
法人北見工業大学内
Fターム(参考) 2C001 AA17 BA06 BB01 BC01 BC08
CB00 CB01 CB02 CB06 CC01
CC08

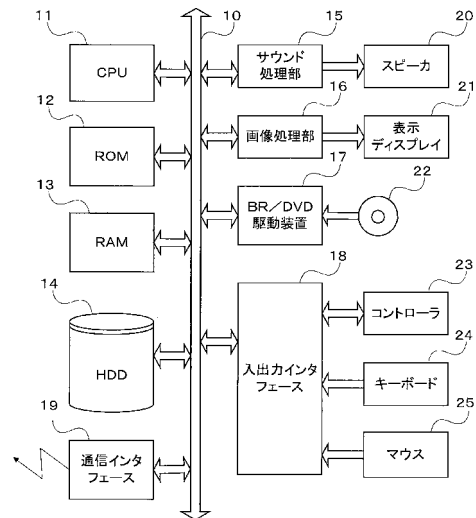
(54) 【発明の名称】 ロールプレイングゲームの攻略法算出装置、算出方法、算出プログラム及びこのプログラムを記録した記録媒体

(57) 【要約】

【課題】被験者の体験データを取得することなく、RPGの攻略法を算出することができ、人間が楽しいと感じるゲーム要素を工学的に把握することができ、コンピュータによって操作されるキャラクタの行動選択の仕方を容易にプログラミングすることができるRPGの攻略法算出装置、算出方法、算出プログラム及びこのプログラムを記録した記録媒体を提供する。

【解決手段】RPGの攻略法算出装置は、RPGのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、この制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出部と、プレイヤーの状態及び時点が与えられるとその時点のその状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定部とを備えており、プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力する。

【選択図】図1



【特許請求の範囲】**【請求項 1】**

ロールプレイングゲームのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、該制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出部と、

前記プレイヤーの状態及び時点が与えられると当該時点の当該状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定部とを備えており、

前記プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力することを特徴とするロールプレイングゲームの攻略法算出装置。

10

【請求項 2】

前記最適政策算出部は、前記制御期間におけるマルコフ決定過程問題を動的計画法で求めるように構成されていることを特徴とする請求項 1 に記載の攻略法算出装置。

【請求項 3】

ロールプレイングゲームのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、該制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出工程と、

前記プレイヤーの状態及び時点が与えられると当該時点の当該状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定工程とを備えており、

前記プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力することを特徴とするロールプレイングゲームの攻略法算出方法。

20

【請求項 4】

前記最適政策算出工程は、前記制御期間におけるマルコフ決定過程問題を動的計画法で求めるように構成されていることを特徴とする請求項 3 に記載の攻略法算出方法。

【請求項 5】

ロールプレイングゲームのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、該制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出手順と、

前記プレイヤーの状態及び時点が与えられると当該時点の当該状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定手順とをコンピュータで実行させ、

前記プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力することを特徴とするロールプレイングゲームの攻略法算出プログラム。

30

【請求項 6】

前記最適政策算出手順は、前記制御期間におけるマルコフ決定過程問題を動的計画法で求めるように構成されていることを特徴とする請求項 5 に記載の攻略法算出プログラム。

40

【請求項 7】

ロールプレイングゲームのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、該制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出手順と、

前記プレイヤーの状態及び時点が与えられると当該時点の当該状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定手順とをコンピュータで実行させる攻略法算出プログラムを記録したコンピュータ読み取り可能な記録媒体であり、

前記プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力することを特徴とするロールプレイングゲー

50

ムの攻略法算出プログラムを記録した記録媒体。

【請求項 8】

前記最適政策算出手順は、前記制御期間におけるマルコフ決定過程問題を動的計画法で求めるように構成されていることを特徴とする請求項 7 に記載の記録媒体。

【発明の詳細な説明】

【技術分野】

【0001】

本発明は、ロールプレイングゲーム（RPG）をモデル化してゲームの攻略法（行動選択の仕方）を算出するための攻略法算出装置、算出方法、算出プログラム及びこのプログラムを記録した記録媒体に関する。

10

【背景技術】

【0002】

近年、コンピュータの低価格化に伴い、テレビゲーム機が広く普及し、ゲームの一分野として、RPGも広く普及している。このようなRPGの開発を行う場合、従来は以下のような問題点を有していた。

（1）プレイヤーが遊ぶ際にどのようにプレイするかを把握するためには、多くの被験者に遊んでもらいプレイヤーの体験データを取得することが行われる。しかしながら、これは、多くの被験者を雇う必要があるため、コストが大幅に高くなるのみならず、データを得るのに多大な時間を要する。

20

（2）日々販売されるRPGは多数あるが、その中でヒットするRPGは非常に少数である。そのようにヒットしたRPGには、人間が楽しいと感じる要素が多数含まれていると考えられるが、どのような要素が楽しいと感じる要素であるのか工学的には未だ把握されていない。

（3）最近では、プレイヤーの補佐を行うキャラクタ（お仲間キャラクタ）をコンピュータが操作するRPGも多いが、このようなお仲間キャラクタの賢い行動（コマンド）選択の仕方をプログラミングすることはかなり難しい。

【0003】

一方、ゲーム情報学等の工学分野においては、RPGをマルコフ決定過程（MDP）等の確率モデルを用いてモデル化し、ゲームを攻略する戦略について数理工学的に扱う研究がなされ、報告されている（例えば、非特許文献 1 及び 2）。

30

【先行技術文献】

【非特許文献】

【0004】

【非特許文献 1】高木幸一郎、雨宮真人、「ロールプレイングゲーム（RPG）の戦闘におけるバランス自動調整システム開発のための基礎的考察」、情報処理学会研究報告、GI、2001（28）、pp. 31 - 38（2001）

【非特許文献 2】高木幸一郎、雨宮真人、「ロールプレイングゲーム（RPG）のバランスとは何か、分析およびその調整に関する提案」、情報処理学会研究報告、GI、2001（58）、pp. 67 - 74（2001）

【発明の概要】

40

【発明が解決しようとする課題】

【0005】

非特許文献 1 及び 2 に開示されているごとき、従来のモデル化は、RPGの一部のみのモデル化であった。即ち、プレイヤーがマップ上を移動するマップモードと、マップモードにおいて敵と遭遇した際に開始される戦闘モードとからなる冒険型のRPGにおいて、戦闘モードのみのモデル化であり、これでは、RPG全体の攻略法（行動選択の仕方）を算出することができず、前述した（1）～（3）のごとき開発支援上の問題点を解決することができなかった。

【0006】

従って本発明の目的は、被験者の体験データを取得することなく、RPGの攻略法（行

50

動選択の仕方)を算出することができるRPGの攻略法算出装置、算出方法、算出プログラム及びこのプログラムを記録した記録媒体を提供することにある。

【0007】

本発明の他の目的は、人間が楽しいと感じるゲーム要素を工学的に把握することができるRPGの攻略法算出装置、算出方法、算出プログラム及びこのプログラムを記録した記録媒体を提供することにある。

【0008】

本発明のさらに目的は、コンピュータによって操作されるキャラクタの行動選択の仕方を容易にプログラミングすることができるRPGの攻略法(行動選択の仕方)を算出することができるRPGの攻略法算出装置、算出方法、算出プログラム及びこのプログラムを記録した記録媒体を提供することにある。

10

【課題を解決するための手段】

【0009】

本発明によれば、RPGのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、この制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出部と、プレイヤーの状態及び時点が与えられるとその時点のその状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定部とを備えており、プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力するRPGの攻略法算出装置が提供される。

20

【0010】

RPGの自動開発に資する研究開発としてMDPなる確率モデルを用いた定式化は従来より行われていたが、従来技術では、RPGの一部のイベントを定式化するのみの不十分なものであった。

【0011】

本発明によれば、RPG全体をMDPで定式化している。これにより、
(A)開発者のみが知っている情報を既知と仮定したもとの攻略法の算出を行い、
(B)開発者のみが知っている情報を未知と仮定した(一般のプレイヤーと同様の立場を仮定した)もとの攻略法の算出を行っている。

30

【0012】

これは、ゲーム情報学分野における学際的な知見を与えるのみならず、ゲーム産業においても以下のごとく有用である。

【0013】

第1に、本発明で算出される目的(お金の獲得等)のもとの攻略法(行動選択の仕方)をコンピュータにシミュレーションさせることによって、その目的におけるプレイヤーのゲーム結果をシミュレートできる。このシミュレーションを利用することによって、マップ上に隠されたアイテムやイベントに遭遇する割合(仮にプレイヤーが1万人いた場合に何人が遭遇できるか)等を把握でき、その割合を見ながら適切な隠し場所を設定することができる。

40

【0014】

第2に、本発明で算出される数理工学的に最適な攻略法と実際のプレイヤーによる攻略法との比較を行うことにより、人間が楽しいと感じるゲーム要素を工学的に把握できる可能性がある。人間が楽しいと感じるゲーム要素を確率モデル上のパラメータ設定のある種のパターンとして把握できれば、そのような要素を多く含むゲーム開発を行うことができる。

【0015】

第3に、近年では、ゲーム中にプレイヤーに協力するコンピュータ操作のキャラクタの登場が多い。このようなコンピュータによって操作されるキャラクタの行動選択の仕方をプログラミングするのは難しかったが、本発明により種々の目的毎の攻略法を算出することによって、各目的に適した(プレイヤーに協力する)キャラクタの行動選択の仕方(攻

50

略法)をプログラミングすることが可能となる。

【0016】

最適政策算出部は、制御期間におけるMDP問題を動的計画法(DP)で求めるように構成されていることが好ましい。

【0017】

本発明によれば、さらに、RPGのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、この制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出工程と、プレイヤーの状態及び時点が与えられるとその時点のその状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定工程とを備えており、プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力するRPGの攻略法算出方法が提供される。

10

【0018】

最適政策算出工程は、制御期間におけるMDP問題をDPで求めるように構成されていることが好ましい。

【0019】

本発明によれば、さらにまた、RPGのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、この制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出手順と、プレイヤーの状態及び時点が与えられるとその時点のその状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定手順とをコンピュータで実行させ、プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力するRPGの攻略法算出プログラムが提供される。

20

【0020】

最適政策算出手順は、制御期間におけるMDP問題をDPで求めるように構成されていることが好ましい。

【0021】

本発明によれば、さらに、RPGのプレイヤーのゲーム開始時点における初期状態と制御期間の長さとは与えられた際に、この制御期間中に得られる総報酬に相当する期待総利得を最大にする政策を出力する最適政策算出手順と、プレイヤーの状態及び時点が与えられるとその時点のその状態においてそれ以降の期待総利得を最大にする最適行動及び期待総利得の最大値を出力する行動決定手順とをコンピュータで実行させる攻略法算出プログラムを記録したコンピュータ読み取り可能な記録媒体であり、プレイヤーの初期状態及び制御期間の長さに対して制御期間における期待総利得を最大にすることが保証された最適政策を出力するRPGの攻略法算出プログラムを記録した記録媒体が提供される。

30

【0022】

最適政策算出手順は、制御期間におけるMDP問題をDPで求めるように構成されていることが好ましい。

【発明の効果】

【0023】

本発明によれば、制御期間における期待総利得を最大にすることが保証された政策が出力されるので、プレイヤーの初期状態と制御期間長とに対して制御期間における期待総利得を最大にする政策を出力することが可能となり、そのRPGに関する攻略法(行動選択の仕方)を算出することができる。

40

【0024】

即ち、本発明によれば、攻略法をコンピュータにシミュレーションさせることによって、被験者の体験データを取得することができ、算出される数理工学的に最適な攻略法と実際のプレイヤーによる攻略法との比較を行うことにより、人間が楽しいと感じるゲーム要素を工学的に把握することができ、さらに、人間が楽しいと感じるゲーム要素を確率モデル上のパラメータ設定のある種のパターンとして把握できれば、そのような要素を多く含

50

むゲーム開発を行うことができる。また、種々の目的毎の攻略法を算出することによって、各目的に適したプレイヤーに協力するキャラクタの攻略法を容易にプログラミングすることが可能となる。

【図面の簡単な説明】

【0025】

【図1】本発明の第1の実施形態として、RPGの開発支援に用いる攻略法算出装置の全体構成を概略的に示すブロック図である。

【図2】図1の実施形態における攻略法算出装置の主要部の動作を説明するフローチャートである。

【図3】図1の実施形態における攻略法算出装置の主要部の構成を概略的に示すブロック図である。

【図4】図1の実施形態におけるDPグラフの一例を示す図である。

【図5】図1の実施形態における攻略法算出装置の行動決定部の動作を説明するフローチャートである。

【図6】第1の実施形態及び第2の実施形態の実施例におけるマップの構成例を示す図である。

【発明を実施するための形態】

【0026】

図1は本発明の第1の実施形態としてRPGの開発支援に用いる攻略法算出装置の全体構成を概略的に示しており、図2は本実施形態における攻略法算出装置の主要部の動作を説明しており、図3は本実施形態における攻略法算出装置の主要部の構成を概略的に示しており、図4は本実施形態におけるDPグラフの一例を示しており、図5は本実施形態における攻略法算出装置の行動決定部の動作を説明している。この第1の実施形態は、各種確率分布を支配する真のパラメータ * が既知の場合である。

【0027】

図1に示すように、本実施形態における攻略法算出装置は、バス10を介して互いに接続された中央処理装置(CPU)11と、リードオンリメモリ(ROM)12と、ランダムアクセスメモリ(RAM)13と、ハードディスク駆動装置(HDD)14と、サウンド処理部15と、画像処理部16と、ブルーレイディスク/デジタルバーサタイルディスク(BR/DVD)駆動装置17と、入出力インタフェース18と、通信インタフェース19とを備えたコンピュータ及びこれを作動させるプログラムから構成される。

【0028】

サウンド処理部15はスピーカ20に接続されており、画像処理部16は表示ディスプレイ21に接続されている。BR/DVD駆動装置17はブルーレイディスク/デジタルバーサタイルディスク/コンパクトディスク(BR/DVD/CD)22が装着可能となっており、入出力インタフェース18にはコントローラ23、キーボード24及びマウス25が接続されている。

【0029】

CPU11は、ROM12に記憶されているオペレーションシステム(OS)やブートプログラム等の基本プログラムに従ってRAM13に記憶されているプログラムを実行して本実施形態の処理を行う。また、CPU11は、RAM13、HDD14、音声処理部15、画像処理部16、BR/DVD駆動装置17、入出力インタフェース18、及び通信インタフェース19の動作を制御する。

【0030】

RAM13は攻略法算出装置のメインメモリとして使用され、HDD14やBR/DVD駆動装置17から転送されたプログラムやデータを記憶する。また、RAM13は、プログラム実行時の各種データが一時的に記憶されるワークエリアとしても使用される。

【0031】

HDD14は、プログラム及びデータがあらかじめ記憶されているか、又は通信インタフェース19を介して外部のネットワーク取り込んだプログラム及びデータが記憶される

10

20

30

40

50

。

【0032】

サウンド処理部15は、CPU11の指示に従ってゲームの背景音や効果音等のサウンドデータを再生するための処理を行い、スピーカ20へその音声信号を出力する。

【0033】

画像処理部16は、CPU11の指示に従って2次元又は3次元グラフィック処理を行い、画像データを生成する。生成された画像データは、表示ディスプレイ21に出力される。表示ディスプレイ21ではなく図示しないTVの画面に表示する場合には、同期信号を付加したビデオ信号を出力する。

【0034】

BR/DVD駆動装置17は、CPU11の指示に従って、セットされたBR/DVD/CD22からゲームに関連するプログラムやデータを読み出し、RAM13へ転送する。また、セットされたBR/DVD/CD22へプログラムやデータの書き込みをすることも可能である。

【0035】

入出力インタフェース18は、コントローラ23、キーボード24及びマウス25とCPU11又はRAM13との間のデータのやり取りを制御する。コントローラ23には、ゲームを行う際に操作される方向キーやボタン等を備えている。

【0036】

通信インタフェース19は、通信回線を介して外部ネットワークに接続されており、CPU11の指示に従って、外部ネットワークとの間でプログラムやデータのやり取りが可能となっている。

【0037】

このような構成の攻略法算出装置において、CPU11は、作動時は、まず、RAM13内にプログラム記憶領域、データ記憶領域及びワークエリアを確保し、HDD14又は外部からプログラム及びデータを取り込んで、プログラム記憶領域及びデータ記憶領域に格納する。次いで、このプログラム記憶領域に格納されたプログラムに基づいて、図2に示す処理を実行する。CPU11がプログラムを実行することによって、図3に概略的に示すとき攻略法算出装置が構築される。

【0038】

即ち、図3に示すように、本実施形態の攻略法算出装置は、最適政策算出部30と、行動決定部31とを備えるように構築される。ここで、最適政策算出部30はDPグラフ作成器30aと、DP実施器30bとを備え、行動決定部31は行動決定器31aと、遷移確率テーブル31bと、利得テーブル31cとを備えるように構築される。

【0039】

この攻略法算出装置の各部説明を行う前に、本実施形態で用いるマルコフ決定過程(MDP)を利用したロールプレイングゲーム(RPG)の数理モデルと、その数理モデルで表現されるRPGの仕様とについて説明する。

【0040】

まず、本実施形態におけるRPGの仕様について説明する。

【0041】

プレイヤーはヒットポイント(HP)と呼ばれる数値を持ち、HPが0となると、次の期にマップ上のスタート位置から再開する。再開時には、HPは、スタート時と同じ最大値 M_{hp} まで回復する。

【0042】

s_{m_i} はマップ上の位置を示し、 SM 、 $SM = \{s_{m_1}, s_{m_2}, \dots, s_{m_{|SM|}}\}$ は、マップ上の位置の集合である。ここで、ゲーム開始時のスタート位置を s_{m_1} とする。なお、本実施形態においては、スタート位置及び現在のプレイヤーの位置は、プレイヤーによって既知であるとする。 f_i は、マップ上の地形の種類を示し、 F 、 $F = \{f_1, f_2, \dots, f_{|F|}\}$ はマップ上の地形の種類の集合である。マップ上の各位置

10

20

30

40

50

がどの地形に該当するかは、関数 $F(s_{m_i})$ F で分かる。

【0043】

e_i は、敵の種類を示し、 E 、 $E = \{e_1, e_2, \dots, e_{|E|}\}$ は敵の種類の集合である。 $M(e_i)$ は、敵 e_i の出現時のこの敵 e_i のHPを示す。プレイヤーは、敵を攻撃することによって敵のHPを0以下にすると、その敵を倒し、その敵に該当する報酬 $G(e_i)$ を得る。

【0044】

プレイヤーが選択できる行動(コマンド)は、マップモードと戦闘モードとでは異なり、マップモードでは a_1 から a_4 が選択可能であり、戦闘モードでは a_5 及び a_6 が選択可能である。 a_1, a_2, a_3, a_4 は、マップ上でそれぞれ右、左、上、下に移動するための行動である。 $mv(s_{m_i}, a_j)$ はプレイヤーが位置 s_{m_i} で行動 a_j を選択した際の移動先位置である。プレイヤーの移動に際して、確率 $p(e_k | F(mv(s_{m_i}, a_j), *))$ で移動先 $mv(s_{m_i}, a_j)$ に敵 e_k が出現し、戦闘モードになる。敵は同時に複数出現することなく、確率 $1 - \sum_{e_k \in E} p(e_k | F(mv(s_{m_i}, a_j), *))$ で何も出現せずにマップモードが続く。

10

【0045】

戦闘モードの行動 a_5 はプレイヤーが戦うための行動であり、確率 $p(C(e_i) | a_5, e_i, *)$ で敵 e_i への攻撃に成功し、その場合、敵 e_i のHPが $C(e_i)$ だけ減少する。プレイヤーは、確率 $1 - p(C(e_i) | a_5, e_i, *)$ で敵 e_i への攻撃に失敗する。また、行動 a_5 の選択とは直接的に関係しないが、戦闘モードでは敵もプレイヤーに対して攻撃し、確率 $p(B(e_i) | e_i, *)$ で敵 e_i がプレイヤーへの攻撃に成功し、その場合、プレイヤーのHPが $B(e_i)$ だけ減少する。攻撃は、プレイヤーが常に先攻すると仮定する。敵 e_i は、確率 $1 - p(B(e_i) | e_i, *)$ でプレイヤーへの攻撃に失敗する。行動 a_6 は、プレイヤーが敵から逃げるための行動であり、確率 $p(map | a_6, *)$ でプレイヤーは次の期にマップモードで移動し、確率 $1 - p(map | a_6, *)$ で戦闘モードが続く。行動 a_6 を選択した場合にも、敵は攻撃してくる。よって、プレイヤーが逃げることに失敗し、かつ敵が攻撃に成功すると、プレイヤーはダメージを受ける。 $*$ は、上述の各確率分布を支配する真のパラメータであり、本実施形態では既知であるとする。

20

【0046】

次に、確率システムの動的な最適化問題を定式化する優れた能力を有する数理モデルであるMDPについてその概要を説明する。

30

【0047】

MDPについては、例えば、金子哲夫、「マルコフ決定理論入門」、槇書店(1973)や森村英典、高橋幸雄、「マルコフ解析」、日科技連、東京(1979)等に記載されている。

【0048】

MDPは、状態 $s_i, s_i \in S$ 、 $S = \{s_1, s_2, \dots, s_{|S|}\}$ ($|S|$ は有限)、各状態で選択できる行動 $a_i, a_i \in A$ 、 $A = \{a_1, a_2, \dots, a_{|A|}\}$ ($|A|$ は有限)、状態 s_i で行動 a_j を選択したもつで、状態 s_k へ遷移する遷移確率 $p(s_k | s_i, a_j, *)$ ($*$ は遷移確率分布を支配する真のパラメータ)、遷移に伴って発生する利得 $r(s_i, a_j, s_k)$ で構成される。MDPの目的は、行動を選び、状態が遷移し、利得を得るという一連のプロセスを繰り返しながら総利得を最大化することである。プロセスの繰り返し回数が有限の場合には、総利得の期待値(期待総利得)を最大化する最適な決定関数を動的計画法(DP)によって求めることができる。具体的には、真のパラメータ $*$ が既知の場合であれば、下記の式(1)を用いて、 t 期の状態が s_i という条件下における t 期以降の期待総利得の最大値 $V(s_i, t)$ を逐次的に計算できる。決定関数は、状態と期とを受け取って、その期で選ぶべき行動を返す関数である。

40

【0049】

50

【数 1】

$$V(s_i, t) = \max_{a_j \in A} \sum_{s_k \in S} p(s_k | s_i, a_j, \xi^*) (r(s_i, a_j, s_k) + V(s_k, t + 1)). \quad (1)$$

【0050】

次に、MDPと本実施形態におけるRPGとの対応について説明する。

【0051】

x_t は、MDPにおけるt期の状態を示す変数であり、式(2)のように構成される。

【0052】

$$x_t = (x_{t,1}, x_{t,2}, x_{t,3}, x_{t,4}) \quad (2)$$

10

ただし、 $x_{t,1}$ はt期におけるプレイヤーのHP、 $x_{t,2}$ はt期におけるプレイヤーのマップ上での位置、 $x_{t,3}$ はt期における敵の種類、 $x_{t,4}$ はt期における敵のHPをそれぞれ示し、マップモードの場合には敵は存在せず、 $x_{t,3} = x_{t,4} = 0$ とする。

【0053】

$A(x_t)$ は、状態 x_t において選択可能なMDPの行動集合を示す。 y_t はMDPにおけるt期に選択した行動を示す変数である。

【0054】

次に、マップモードのt期の状態 x_t で行動 y_t を選択したときの状態遷移について説明する。t+1期には、確率 $p(e_i | F(mv(x_{t,2}, y_t), \xi^*))$ で敵 e_i 20
が出現し、戦闘モードの状態 x_{t+1} 、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, mv(x_{t,2}, y_t), e_i, M(e_i)) \end{aligned} \quad (3)$$

に遷移する。ただし、ゲームのスタート位置である sm_1 が、移動先 $mv(x_{t,2}, y_t)$ の場合には敵は出現しない($p(e_i | F(mv(x_{t,2}, y_t), \xi^*)) = 0$)とする。また、確率 $1 - p(e_i | F(mv(x_{t,2}, y_t), \xi^*))$ で敵が出現せずにマップモードの状態 x_{t+1} に遷移する。このときの状態 x_{t+1} は、移動先 $mv(x_{t,2}, y_t)$ が sm_1 の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (M_{hp}, sm_1, x_{t,3}, x_{t,4}) \end{aligned} \quad (4)$$

30

移動先 $mv(x_{t,2}, y_t)$ が sm_1 以外の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, mv(x_{t,2}, y_t), x_{t,3}, x_{t,4}) \end{aligned} \quad (5)$$

である。式(4)の場合は、プレイヤーがスタート位置 sm_1 に戻り、HPを最大値 M_{hp} まで回復した状態である。

【0055】

次に、戦闘モードのt期の状態 x_t で行動 y_t を選択したときの状態遷移について、行動 y_t が行動 a_5 (戦う)の場合と行動 a_6 (逃げる)の場合とに分けて説明する。

【0056】

まず、行動 a_5 (戦う)の場合について説明する。確率 $1 - p(C(x_{t,3}) | a_5, x_{t,3}, \xi^*) (1 - p(B(x_{t,3}) | x_{t,3}, \xi^*))$ でプレイヤーと敵との両方が攻撃に失敗し、状態 x_{t+1} 、

40

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, x_{t,2}, x_{t,3}, x_{t,4}) \end{aligned} \quad (6)$$

に遷移する。確率 $(1 - p(C(x_{t,3}) | a_5, x_{t,3}, \xi^*)) p(B(x_{t,3}) | x_{t,3}, \xi^*)$ でプレイヤーは攻撃に失敗し、敵は攻撃に成功し、状態 x_{t+1} へ遷移する。

このときの状態 x_{t+1} は、 $x_{t,1} > B(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1} - B(x_{t,3}), x_{t,2}, x_{t,3}, x_{t,4}) \end{aligned} \quad (7)$$

50

で、 $x_{t,1} \leq B(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (M_{hp}, sm_1, 0, 0) \end{aligned} \quad (8)$$

である。式(8)の場合には、プレイヤーが敵に倒されて、ゲームのスタート位置 sm_1 からの再開である。確率 $p(C(x_{t,3}) | a_5, x_{t,3}, *) (1 - p(B(x_{t,3}) | x_{t,3}, *))$ でプレイヤーは攻撃に成功し、敵は攻撃に失敗し、状態 x_{t+1} へ遷移する。このときの状態 x_{t+1} は、 $x_{t,4} > C(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, x_{t,2}, x_{t,3}, x_{t,4} - C(x_{t,3})) \end{aligned} \quad (9)$$

で、 $x_{t,4} \leq C(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, x_{t,2}, 0, 0) \end{aligned} \quad (10)$$

である。式(10)の場合には、敵 $x_{t,3}$ を倒すことに成功しているので、この状態遷移に伴い、利得 $r(x_t, a_5, x_{t+1}) = G(x_{t,3})$ を得る。確率 $p(C(x_{t,3}) | a_5, x_{t,3}, *) p(B(x_{t,3}) | x_{t,3}, *)$ でプレイヤーと敵の両方が攻撃に成功し、状態 x_{t+1} へ遷移する。このときの状態 x_{t+1} は、 $x_{t,4} > C(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, x_{t,2}, 0, 0) \end{aligned} \quad (11)$$

で、 $x_{t,4} > C(x_{t,3})$ かつ $x_{t,1} > B(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1} - B(x_{t,3}), x_{t,2}, x_{t,3}, x_{t,4} - C(x_{t,3})) \end{aligned} \quad (12)$$

で、 $x_{t,4} > C(x_{t,3})$ かつ $x_{t,1} \leq B(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (M_{hp}, sm_1, 0, 0) \end{aligned} \quad (13)$$

である。式(11)の場合には、敵 $x_{t,3}$ を倒すことに成功しているので、この状態遷移に伴い、利得 $r(x_t, a_5, x_{t+1}) = G(x_{t,3})$ を得る。式(13)の場合にはプレイヤーが敵に倒されて、ゲームのスタート位置 sm_1 からの再開である。

【0057】

次に、戦闘モードの t 期の状態 x_t で行動 a_6 (逃げる) を選択したときの状態遷移について説明する。確率 $p(map | a_6, *)$ でプレイヤーが逃げることに成功し、状態 x_{t+1} 、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, x_{t,2}, 0, 0) \end{aligned} \quad (14)$$

に遷移する。確率 $1 - p(map | a_6, *) (1 - p(B(x_{t,3}) | x_{t,3}, *))$ でプレイヤーが逃げることに失敗し、敵が攻撃に失敗し、状態 x_{t+1} 、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1}, x_{t,2}, x_{t,3}, x_{t,4}) \end{aligned} \quad (15)$$

に遷移する。確率 $(1 - p(map | a_6, *)) p(B(x_{t,3}) | x_{t,3}, *)$ でプレイヤーが逃げることに失敗し、敵が攻撃に成功し、状態 x_{t+1} へ遷移する。このときの状態 x_{t+1} は、 $x_{t,1} > B(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (x_{t,1} - B(x_{t,3}), x_{t,2}, x_{t,3}, x_{t,4}) \end{aligned} \quad (16)$$

で、 $x_{t,1} \leq B(x_{t,3})$ の場合には、

$$\begin{aligned} x_{t+1} &= (x_{t+1,1}, x_{t+1,2}, x_{t+1,3}, x_{t+1,4}) \\ &= (M_{hp}, sm_1, 0, 0) \end{aligned} \quad (17)$$

である。式(17)の場合には、プレイヤーが敵に倒されて、ゲームのスタート位置 sm_1 からの再開である。

10

20

30

40

50

【 0 0 5 8 】

前述した通り、プレイヤーが敵 $x_{t, 3}$ を倒した状態遷移に伴う利得は、 $r(x_t, a_5, x_{t+1}) = G(x_t, 3)$ である。その他の状態遷移に伴う利得は、 $r(x_t, a_5, x_{t+1}) = 0$ である。本実施形態では、初期状態 x_1 が $x_1 = (M_{hp}, S_{m1}, 0, 0)$ であり、各期の状態は観測可能である。また、プレイヤーや敵の攻撃力 $C(e_i)$ 、 $B(e_i)$ 及び敵を倒したときの報酬 $G(e_i)$ 等は全て既知であるとする。このもとで、T 期間のプレイを行って総利得

【 数 2 】

$$\sum_{t=1}^T r(x_t, y_t, x_{t+1})$$

10

の最大化を目的とする。

【 0 0 5 9 】

次に、図 2 に示されたフローチャート、図 3 に示されたブロック図、図 4 に示された図、及び図 5 に示されたフローチャートを参照して、本実施形態の攻略法算出装置の動作を説明する。

【 0 0 6 0 】

図 2 及び図 3 に示すように、まず、構築された最適政策算出部 30 における DP グラフ作成器 30 a にプレイヤーの初期状態 x_1 と制御期間長 T とが入力される (ステップ S 1)。この初期状態 x_1 と制御期間長 T とは HDD 14 に格納されているデフォルト値又は前回のプレイ結果値を RAM 13 のデータ記憶領域に格納したものであっても良いし、キーボード 24 から入力した値であっても良い。

20

【 0 0 6 1 】

プレイヤーの初期状態 x_1 と制御期間長 T とが入力されると、DP グラフ作成器 30 a は、T 期間の期待総利得を最大化するための動的計画法 (DP) の問題を解くための DP グラフを作成する (ステップ S 2)。例えば、想定されるプレイヤーの全状態を要素とする状態集合 S が、 $S = \{s_1, s_2, s_3, s_4\}$ で $x_1 = s_1$ の場合であれば、図 4 のような DP グラフが作成される。これは、1 時点目 (1 期) はプレイヤーの初期状態で表現され、2 時点目から T 時点目まではプレイヤーの想定される各状態で表現されたグラフにおいて、末端の T 時点目 (T 期) のノードから遡りながら DP で T 期間の MDP 問題を解くことによって、T 期間の期待総利得を最大化する最適政策を求めるための準備である。

30

【 0 0 6 2 】

次いで、DP 実施器 30 b が DP によって T 期間の MDP 問題を解くことによって、T 期間の期待総利得を最大化する最適政策が求められる (ステップ S 3)。DP 実施器 30 b は、DP グラフの末端の各ノードから順にそのノードでの最適な行動 (RPG におけるプレイヤーのコマンド選択) とそのノード以降の期待総利得の最大値を、行動決定部 31 における行動決定器 31 a と連携して求める (ステップ S 4)。

【 0 0 6 3 】

即ち、各ノード毎にそのノードの時点 t (何時点目かを示す自然数) とプレイヤーの状態 x_t (t 時点目のプレイヤーの状態) とを行動決定器 31 a へ送ると、そのノードにおける最適な行動とそのノード以降の期待総利得の最大値とが求められて行動決定器 31 a から送り返される。

40

【 0 0 6 4 】

その後、DP グラフの 1 時点目のノードまで全て解き終わったかが判断され (ステップ S 5)、解き終わっていれば DP グラフの全ノードにおける最適な行動とそのノード以降の期待総利得の最大値が最適政策として出力される (ステップ S 6)。

【 0 0 6 5 】

次に、図 3 及び図 5 に示されたフローチャートを参照して行動決定部 31 の動作、即ち図 2 におけるステップ S 4 の動作を説明する。

【 0 0 6 6 】

50

まず、最適政策算出部 30 の DP 実施器 30 b から行動決定部 31 の行動決定器 31 a へ、時点 t (何時点目かを示す自然数) とプレイヤーの状態 x_t (t 時点目のプレイヤーの状態) とが入力される (ステップ S 41)。

【0067】

次いで、入力された時点 t とプレイヤーの状態 x_t とに応じてそのノード以降の期待総利得の最大値とそのノードにおける最適な行動とが算出される (ステップ S 42)。 $t = T$ の場合には、式 (18) で、そのノード以降の期待総利得の最大値が求められる。

【0068】

【数 3】

$$V(x_T, T) = \max_{y_T \in A(x_T)} \sum_{x_{T+1}} p(x_{T+1} | x_T, y_T, \theta^*) r(x_T, y_T, x_{T+1}), \quad (18)$$

10

ただし、 $V(x_T, T)$ は T 時点目の状態 x_T から時点 $T+1$ への状態遷移に伴う最後の 1 期間の期待総利得の最大値である。 $p(x_{T+1} | x_T, y_T, \theta^*)$ は遷移確率テーブル 31 b から読み取ったものである。 $r(x_T, y_T, x_{T+1})$ は利得テーブル 31 c から読み取ったものである。 $1 \leq t \leq T-1$ の場合には次の式 (19) でそのノード以降の期待総利得の最大値が求められる。

【0069】

【数 4】

$$V(x_t, t) = \max_{y_t \in A(x_t)} \sum_{x_{t+1}} p(x_{t+1} | x_t, y_t, \theta^*) (r(x_t, y_t, x_{t+1}) + V(x_{t+1}, t+1)), \quad (19)$$

20

ただし、 $V(x_T, T)$ は t 時点目の状態が x_t という条件のもとでの、 t 時点以降の期待総利得の最大値である。本実施形態では DP を利用しているので、このように部分最適解を再利用している。 $t = T$ の場合には次の式 (20) でそのノードにおける最適な行動が求められる。

【0070】

【数 5】

$$d^*(x_T, T) = \arg \max_{y_T \in A(x_T)} \sum_{x_{T+1}} p(x_{T+1} | x_T, y_T, \theta^*) r(x_T, y_T, x_{T+1}), \quad (20)$$

30

ただし、 $d^*(x_T, T)$ は T 時点目の状態 x_T において選択すべき最適な行動である。 $1 \leq t \leq T-1$ の場合には次の式 (21) でそのノードにおける最適な行動が求められる。

【0071】

【数 6】

$$d^*(x_t, t) = \arg \max_{y_t \in A(x_t)} \sum_{x_{t+1}} p(x_{t+1} | x_t, y_t, \theta^*) (r(x_t, y_t, x_{t+1}) + V(x_{t+1}, t+1)), \quad (21)$$

ただし、 $d^*(x_t, t)$ は t 時点目の状態 x_t において選択すべき最適な行動である。

40

【0072】

その後、そのノードにおける最適な行動とそのノード以降の期待総利得の最大値が最適政策算出部 30 の DP 実施器 30 b へ出力され (ステップ S 43)、前述の最適政策が出力されるのである。

【0073】

以上説明したように、第 1 の実施形態によれば、DP を用いて各時点のプレイヤーの各状態において、その時点以降の期待総利得を最大化し、最終的に制御期間における期待総利得を最大にすることが保証された政策が出力されるので、プレイヤーの初期状態と制御期間長とに対して制御期間における期待総利得を最大にする政策を出力することが可能となり、その RPG に関する攻略法 (行動選択の仕方) を算出することができる。

50

【0074】

即ち、本実施形態によれば、攻略法をコンピュータにシミュレーションさせることによって、被験者の体験データを取得することができる。例えば、プレイヤーのゲーム結果をシミュレーションすることによって、マップ上に隠されたアイテムやイベントに遭遇する割合等を把握でき、その割合を見ながら適切な隠し場所を設定する等のゲーム開発支援を行うことができる。また、算出される数理工学的に最適な攻略法と実際のプレイヤーによる攻略法との比較を行うことにより、人間が楽しいと感じるゲーム要素を工学的に把握することができる。人間が楽しいと感じるゲーム要素を確率モデル上のパラメータ設定のある種のパターンとして把握できれば、そのような要素を多く含むゲーム開発を行うことができる。また、種々の目的毎の攻略法を算出することによって、各目的に適したプレイヤーに協力するキャラクタの攻略法を容易にプログラミングすることが可能となる。

10

【0075】

次に、各種確率分布を支配する真のパラメータ θ^* が未知である、本発明を拡張した第2の実施形態について説明する。

【0076】

本実施形態における攻略法算出装置の構成は基本的には、第1の実施形態の場合と同様であり、従ってその構成の説明は省略する。

【0077】

真のパラメータ未知の場合を説明するために、いくつかの新たな定義を行う。 $p(\cdot)$ はパラメータ θ の事前分布であり、既知であるとする。 Θ はパラメータ空間であり、 $\theta^* \in \Theta$ である。 x^t, y^t は t 期目の状態 x_t に至るまでの遷移系列であり、 $x^t, y^t = x_1, y_1, \dots, x_t, y_t$ である。

20

【0078】

真のパラメータ既知の場合には、DPでT時点から遡りながら各時点の各状態に対して行動選択を行うが、真のパラメータ未知の場合には、DPでT時点から遡りながら各時点の各状態と1時点からその時点に至るまでの各遷移系列の組に対して行動選択を行う。

【0079】

T時点目の状態 x_T (全ての状態の候補) と遷移系列 x^T, y^T (全ての遷移系列の候補) の組に対する処理は以下の通りである。

【0080】

【数7】

$$d_B^*(x_T, x^T, y^T, T) = \arg \max_{y_T \in A(x_T)} \sum_{x_{T+1}} \int_{\Theta} p(\theta | x^T, y^T) p(x_{T+1} | x_T, y_T, \theta) d\theta r(x_T, y_T, x_{T+1}), \quad (22)$$

30

$$V_B(x_T, x^T, y^T, T) = \max_{y_T \in A(x_T)} \sum_{x_{T+1}} \int_{\Theta} p(\theta | x^T, y^T) p(x_{T+1} | x_T, y_T, \theta) d\theta r(x_T, y_T, x_{T+1}), \quad (23)$$

40

ただし、 $p(\cdot | x^T, y^T)$ は1時点からT時点に遷移系列 x^T, y^T のように遷移した場合の事後分布である。

【0081】

t時点目 (1 ≤ t ≤ T-1) の状態 x_t (全ての状態の候補) と遷移系列 x^t, y^t (全ての遷移系列の候補) の組に対する処理は以下の通りである。

【0082】

【数 8】

$$d_B^*(x_t, x^t y^{t-1}, t) = \arg \max_{y_t \in A(x_t)} \sum_{x_{t+1}} \int_{\Theta} p(\theta | x^t y^{t-1}) p(x_{t+1} | x_t, y_t, \theta) d\theta \\ (r(x_t, y_t, x_{t+1}) + V_B(x_{t+1}, x^{t+1} y^t, t+1)). \quad (24)$$

$$V_B(x_t, x^t y^{t-1}, t) = \max_{y_t \in A(x_t)} \sum_{x_{t+1}} \int_{\Theta} p(\theta | x^t y^{t-1}) p(x_{t+1} | x_t, y_t, \theta) d\theta \\ (r(x_t, y_t, x_{t+1}) + V_B(x_{t+1}, x^{t+1} y^t, t+1)). \quad (25)$$

10

式(22)から式(25)を用いて $d_B^*(x_1, x_1, 1)$ まで求めることによって、1時点目からT時点目までの全ての状態と遷移系列の組とに対して、ベイズ基準のもとで総利得を最大にするという点で最適な行動選択の仕方を求めることができる。

【0083】

式(22)から式(25)には積分計算が含まれており、一般的に、積分計算は計算量が多いが、二項分布(敵の出現以外の確率分布)の事前分布としてベータ分布を、多項分布(敵の出現の確率分布)の事前分布としてディリクレ分布をそれぞれ仮定すると、積分計算は四則演算で実施することができる(Matsushima, T., Hirasawa, S., A Bayes coding algorithm for Markov models, TECHNICAL REPORT OF IEICE, IT95-1, pp. 1-6 (1995))。四則演算の一例として、マップモードのt期の状態 x_t において行動 y_t を選択したもとの、敵 e_i が出現し、戦闘モードの状態 x_{t+1} に遷移する場合の

20

【数 9】

$$\int_{\Theta} p(\theta | x^t y^{t-1}) p(x_{t+1} | x_t, y_t, \theta) d\theta$$

の計算を以下に示す。

【0084】

【数 10】

$$\int_{\Theta} p(\theta | x^t y^{t-1}) p(x_{t+1} | x_t, y_t, \theta) d\theta = \int_{\Theta} p(\theta | x^t y^{t-1}) p(e_i | F(mv(x_{t,2}, y_t)), \theta) d\theta \\ = \frac{N(F(mv(x_{t,2}, y_t)), e_i | x^t y^{t-1}) + \alpha_1}{N(F(mv(x_{t,2}, y_t)) | x^t y^{t-1}) + \alpha_2}, (26)$$

30

ただし、

$$\alpha_1 = \alpha(e_i | F(mv(x_{t,2}, y_t))), \quad (27)$$

$$\alpha_2 = \sum_{e_j \in E} \alpha(e_j | F(mv(x_{t,2}, y_t))) + \alpha(mv(x_{t,2}, y_t) | F(mv(x_{t,2}, y_t))), \quad (28)$$

40

ここで、 $N(F(mv(x_{t,2}, y_t)), e_i | x^t y^{t-1})$ は系列 $x^t y^{t-1}$ 中で地形の種類が $F(mv(x_{t,2}, y_t))$ の位置で敵 e_i が出現した回数、 $(e_i | F(mv(x_{t,2}, y_t)))$ は $p(e_i | F(mv(x_{t,2}, y_t)))$ に対するディリクレ分布(事前分布)のパラメータ、 $(mv(x_{t,2}, y_t) | F(mv(x_{t,2}, y_t)))$ は

【数 11】

$$1 - \sum_{e_i \in E} p(e_i | F(mv(x_{t,2}, y_t)), \theta)$$

に対するディリクレ分布(事前分布)のパラメータを示す。このように、事前分布としてディリクレ分布やベータ分布を採用することにより、積分計算を四則演算で置き換えるこ

50

とができる。ディリクレ分布やベータ分布のパラメータの設定が事前分布の設定に相当するが、事前に何も情報が無い場合の設定の仕方についてはベイズ統計学やその応用分野で種々の方法が研究されている。多くの分野で良好な性質が報告されているジェフリーズの事前分布が有名であり、本実施形態に適用する場合は、各パラメータを0.5に設定することに相当する。例えば、Berger, J. O., *Statistical Decision Theory and Bayesian Analysis*, Springer-Verlag, New York (1980)、繁榊算男, ベイズ統計入門, 東京大学出版会 (1985)、Matsushima, T., Hirasawa, S., *A Bayesian coding algorithm for Markov models*, TECHNICAL REPORT OF IEICE, IT95-1, pp. 1-6 (1995)、鈴木謙, ベイシアネットワーク入門, 培風館, 東京 (2009) を参照。事前分布にディリクレ分布やベータ分布を採用してジェフリーズの前分布に設定し、式(22)から式(25)で処理することにより、真のパラメータ未知の場合にベイズ基準のもとで総利得を最大化することができる。

【0085】

以上の説明では、真のパラメータ未知の場合のベイズ最適な行動選択の仕方を求めるアルゴリズムについて述べた。事前分布としてディリクレ分布やベータ分布を採用することにより、積分計算を四則演算に置き換えた。しかしながら、ベイズ最適な行動選択の仕方を求めるためには、多大な計算量が必要となる。真のパラメータ既知の場合には、DPの各時点毎に式(19)の処理を状態数分だけ実施すればよい。一方、真のパラメータ未知のベイズ最適の場合には、DPの各時点毎に式(25)の処理を状態数と遷移系列の個数との積分だけ実施する必要があり、処理の回数は時点の数(t期のt)に対する指数オーダーとなる。

【0086】

そこで、本実施形態の望ましい態様として、真のパラメータ未知の場合に近似を行う例を説明する。ここで、学習データLを新たに導入する。学習データは過去のゲームのプレイデータである遷移系列の集合であったり、敵の出現確率などの個々の確率分布について真のパラメータの分布から発生させたサンプルデータであったり、種々の形態が適用可能である。

【0087】

前述したベイズ最適な方法では、各時点毎にその時点までの遷移系列 $x^t y^{t-1}$ に対する事後分布によって、

【数12】

$$\int_{\Theta} p(\theta | x^t y^{t-1}) p(x_{t+1} | x_t, y_t, \theta) d\theta$$

を計算したが、近似アルゴリズムでは、時点に関係なく学習データLによる事後分布を用いて

【数13】

$$\int_{\Theta} p(\theta | L) p(x_{t+1} | x_t, y_t, \theta) d\theta$$

を計算する。具体的には、

【数14】

$$\int_{\Theta} p(\theta | L) p(x_{t+1} | x_t, y_t, \theta) d\theta$$

を真のパラメータ既知の場合の式(18)から式(21)に代入して行動選択の仕方を求める。

【0088】

近似アルゴリズムにより、DPの処理の回数は真のパラメータ既知の場合と同じ回数に軽減することができる。有限の学習データに対する理論的な精度保証はないが、漸近的には学習データによる事後分布を用いた推定値が真のパラメータに収束するので、求める行

10

20

30

40

50

動選択の仕方も真のパラメータ既知の場合に収束する。

【実施例】

【0089】

真のパラメータ θ^* が既知の場合の行動選択の仕方を求めるアルゴリズムに関する第1の実施形態について実験を行った実施例を説明する。図6にこの実施例におけるマップを示す。実験結果が理解しやすいように9マスからなる小規模のマップとした。以下の表1～表5は地形の設定、確率の設定及びその他の設定を示している。

【0090】

【表1】

表1 地形の設定

$F(s m_1)$	$F(s m_2)$	$F(s m_3)$	$F(s m_4)$	$F(s m_5)$
f_1	f_2	f_3	f_2	f_2

$F(s m_6)$	$F(s m_7)$	$F(s m_8)$	$F(s m_9)$
f_3	f_3	f_3	f_3

10

【0091】

【表2】

表2 確率の設定 (その1)

$p(e_1 f_2, \theta^*)$	$p(e_2 f_2, \theta^*)$	$p(e_1 f_3, \theta^*)$	$p(e_2 f_3, \theta^*)$
0.3	0.0	0.0	0.8

20

【0092】

【表3】

表3 確率の設定 (その2)

$p(C(e_1) a_5, e_1, \theta^*)$	$p(C(e_2) a_5, e_2, \theta^*)$	$p(B(e_1) e_1, \theta^*)$
0.9	0.9	0.6

$p(B(e_2) e_2, \theta^*)$	$p(\text{map} a_6, \theta^*)$
0.6	0.6

30

【0093】

【表4】

表4 その他の設定 (その1)

M_{hp}	E	$M(e_1)$	$M(e_2)$	$G(e_1)$	$G(e_2)$
10	$\{e_1, e_2\}$	2	4	1	5

【0094】

【表5】

表5 その他の設定 (その2)

$C(e_1)$	$C(e_2)$	$B(e_1)$	$B(e_2)$	T
3	3	1	3	10

40

【0095】

以上の設定で、真のパラメータ既知の場合のアルゴリズムを適用して、10期間の期待総利得を最大化するための各期における行動選択の仕方を求めた。

【0096】

50

結果の一部について述べると、例えば、時点3(3期)でプレイヤーのHPが6であり位置 s_{m_4} にいるマップモードの状態では、最適の行動選択は a_3 という上のマス s_{m_7} への移動であった。これは、プレイヤーのHPにまだ余裕があるので、弱くて報酬の小さい敵が出現する右(a_1)やHPを回復する下(a_4)ではなく、強くて報酬の大きい敵が出現する上(a_3)への移動を選択しているのである。他方、同じ時点3(3期)のマップモードであってもHPが1であり位置 s_{m_4} にいる状態では、行動 a_4 を選択して下のHPを回復してくれるスタート位置 s_{m_1} へ移動した。これは、プレイヤーのHPに余裕がないので回復するための行動選択である。また、時点9(9期)にHPが1であり位置 s_{m_4} にいるマップモードの状態では、行動 a_1 を選択して、弱くて報酬の小さい敵が出現する右の s_{m_5} へ移動した。これは、プレイヤーのHPに余裕はないが、残りの期間にも余裕がないため、弱い敵が出現する s_{m_5} への移動を選択しているのである。

10

【0097】

このように、真のパラメータ既知の場合のアルゴリズムを適用することにより、開発者であれば知っている真のパラメータの情報を利用して対象期間の期待総利得を最大にする行動選択の仕方を求めることができる。

【0098】

次に、真のパラメータが未知の場合の行動選択の仕方を求めるアルゴリズムに関する第2の実施形態について実験を行った実施例を説明する。真のパラメータ既知の場合の実施例と同じ設定のもとで、真のパラメータ未知の場合の近似アルゴリズムを適用した。学習データとして、各確率分布毎にサンプルデータを発生させた。そのもとで、近似アルゴリズムを適用し、各時点の各状態における行動選択が真のパラメータ既知の場合と比較して一致するかどうか調べた。各確率分布の学習データ数を10、100、1000と変化させ、それぞれの学習データ数に対して100パターンの学習データを発生させて適用実験を行った。真のパラメータ既知の場合の行動選択との一致率は100パターンの平均で、学習データ数が10の場合で約88%、学習データ数が100の場合で約94%、学習データ数が1000の場合で約96%であった。

20

【0099】

このように、学習データによる事後分布を利用する近似アルゴリズムを用いることにより、真のパラメータが未知の場合でも、DPに必要な計算量を、真のパラメータ既知の場合の計算量と同程度とすることができる。また、少ない実験例ではあるが、学習データの増加に伴い真のパラメータ既知の場合との行動選択の一致率が高くなることが確認できた。

30

【0100】

以上述べた実施形態は全て本発明を例示的に示すものであって限定的に示すものではなく、本発明は他の種々の変形態様及び変更態様で実施することができる。従って本発明の範囲は特許請求の範囲及びその均等範囲によってのみ規定されるものである。

【符号の説明】

【0101】

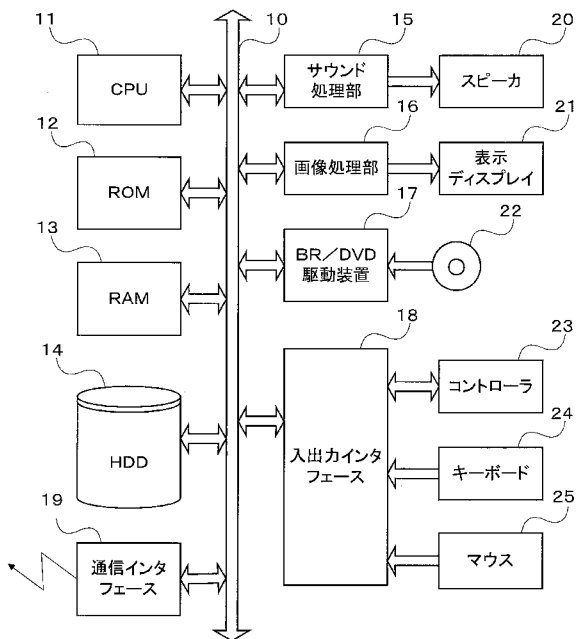
- 10 バス
- 11 CPU
- 12 ROM
- 13 RAM
- 14 HDD
- 15 サウンド処理部
- 16 画像処理部
- 17 BR/DVD駆動装置
- 18 入出力インタフェース
- 19 通信インタフェース
- 20 スピーカ
- 21 表示ディスプレイ

40

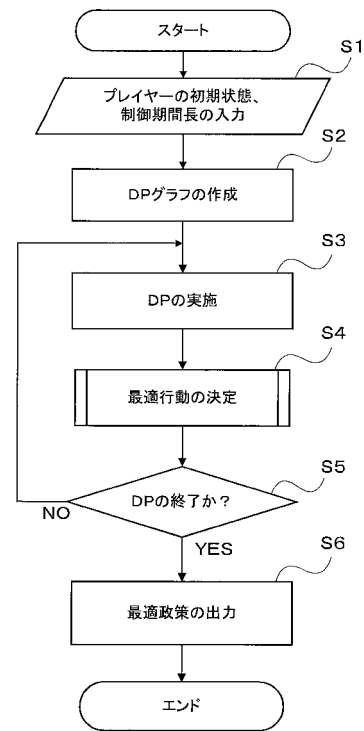
50

- 2 2 BR / DVD / CD
- 2 3 コントローラ
- 2 4 キーボード
- 2 5 マウス
- 3 0 最適政策算出部
- 3 0 a DPグラフ作成器
- 3 0 b DP実施器
- 3 1 行動決定部
- 3 1 a 行動決定器
- 3 1 b 遷移確率テーブル
- 3 1 c 利得テーブル

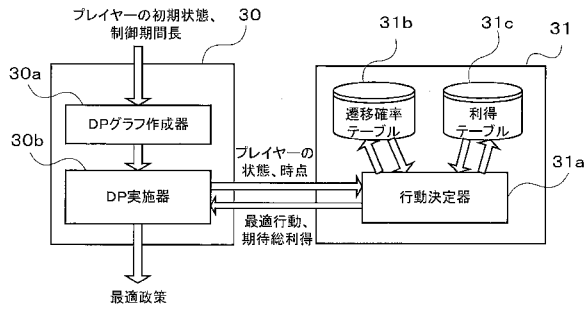
【 図 1 】



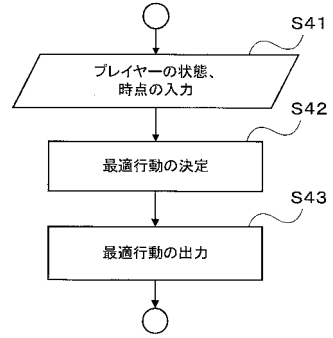
【 図 2 】



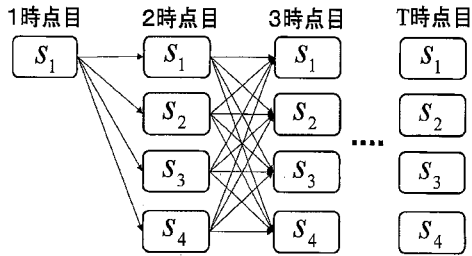
【 図 3 】



【 図 5 】



【 図 4 】



【 図 6 】

sm_7	sm_8	sm_9
sm_4	sm_5	sm_6
sm_1	sm_2	sm_3