

遺伝アルゴリズムによる制約付きマルコフ決定過程の解法

平山 克己*・河合 一

*社会開発工学専攻・社会開発システム工学科

(1995年8月29日受理)

A Solving Method of a MDP with Constraint by Genetic Algorithm

by

Katsumi HIRAYAMA¹⁾・Hajime KAWAI²⁾

¹⁾Course in Engineering of Social Development

²⁾Department of Social Systems Engineering

(Received August 29, 1995)

We consider discrete time Markov decision process (MDP) with finite state space, finite action space and two kinds of immediate rewards. The problem is to maximize time average reward generated by on reward stream, subject to that the other reward is not smaller than a prescribed value. The problem is analyzed in the range of pure stationary policies. MDP with one optimality criterion and no constraint can be solved by usual policy improvement method. MDP with one reward constraint can be solved by linear programming, in the range of mixed policies. On the other hand, however, when we restrict the policies to pure policies the problem is some combinatorial problem, for which any solving method has not been discovered. In this paper, we propose an approach applying Genetic Algorithm in order to carry on a search process effectively and to obtain a near optimal pure stationary policy. A numerical example is given to examine the efficiency of the approach proposed here.

1 はじめに

本論文では、有限状態空間、有限決定空間、及び2種類の直接利得を持つ離散時間マルコフ決定過程 (Markov Decision Process: 略してMDP) を取り扱い、一方の利得から生じる時間平均利得をある与えられた値以上に保護する純定常政策の中で、他方の利得から生じる時間平均利得を最大にする政策を定める制約付MDP問題について考える。

一制約を持つMDPは、既に Beulter and Ross¹⁾により混合戦略の範囲で考察され、最適政策は、せいぜい二つの純政策の混合政策により与えられることが示されている。ただし、混合政策の下では各決定を每期確率的に選択することになり、管理上面倒な点が多く、純政策の範囲内で最適政策を求めることは現実的な意味で重要な問題であると思われる。しかし、純政策に限定すると、組合せの問題となり、厳密解の導出が非常に困難となる。

そこで、本研究では制約付マルコフ決定過程に対して、遺伝アルゴリズムに政策改善法を加味した新しいアプローチを提案する。

2 制約つきMDP

はじめに以下の記号を定義する。

- $I = \{0, 1, \dots, N\}$: 状態空間
- $D_i = \{1, 2, \dots, K_i\}$: 状態*i*における状態空間
- q_{ij}^k : 状態*i*で決定*k*を選択したときの推移確率
- a_{ij}^k, b_{ij}^k : 状態*i*で決定*k*を選択したときに生じる利得
- S : 純政策の集合, すなわち
- $S = D_1 \times D_2 \times \dots \times D_N$
- s : 純政策, すなわち, $s \in S$

ここですべての純政策に対し、マルコフ決定過程は完全エルゴディックであるとす。すなわち、定常分布を持つ。

- $g(s)$: 政策*s*を採用したときの利得 a_{ij}^k から生じる時間平均利得
- $h(s)$: 政策*s*を採用したときの利得 b_{ij}^k から生じる時間平均利得
- π_i^s : 政策*s*を採用したときの定常分布

なお、表現の簡潔化のため、 $q_{ij}^k, a_{ij}^k, b_{ij}^k$ をそれぞれ、政策*s*採用したときの推移確率および、状態*i*における利得を表すとする。

$g(s), h(s)$ はそれぞれ、 $g(s), v_i(s)$ および $h(s), w_i(s)$ ($i \in I$)を未知数とする次の連立方程式

$$\begin{cases} g(s) + v_i(s) = a_i^s + \sum_j q_{ij}^s v_j(s) & i = 1, \dots, N \\ v_0(s) = 0 \end{cases} \quad (1)$$

$$\begin{cases} h(s) + w_i(s) = b_i^s + \sum_j q_{ij}^s w_j(s) & i = 1, \dots, N \\ v_0(s) = 0 \end{cases} \quad (2)$$

の解として与えられる。あるいは、定常分布を用い

$$g(s) = \sum_i \pi_i^s a_i^s \quad (3)$$

$$h(s) = \sum_i \pi_i^s b_i^s \quad (4)$$

$$\pi_j = \sum_i \pi_i^s q_{ij}^s, \quad j \in I \quad (5)$$

$$\sum_i \pi_i = 1 \quad (6)$$

で与えられる。

以上の記号を用いると我々の問題は

$$s \in \{s | h(s) \geq \alpha, s \in S\} \quad g(s) \quad (7)$$

で表現される。

2.1 混合政策と純政策

本研究では触れていないが、例えば図1のように制約付きMDPを混合政策の範囲で考えると、理論的に厳密解が得られることが示されている²⁾。しかし、混合政策は決定を確定的には選ばず、確率的に選ぶ政策である。したがって、意思決定者にとっては純政策の範囲で考える方が現実的であり、取り扱い易いと考えられる。

また、図1のように混合政策では端点を結ぶ直線と時間平均利得*h*の制約値 α が交わる点が最適解となる。しかし、純政策に限定すると実行可能解は離散的な点上に存在し、時間平均利得*h*の制約値 α によっては最適解は端点を結ぶ直線上にあるとは限らず、組合せの問題となっている。そのため、理論的に厳密解を得る手段が現在では存在しない。

3 遺伝アルゴリズムの概要

遺伝アルゴリズムは、1960年代にアメリカのホーランドによって基本的な考え方が提唱された。自然界の生物集団の中で、長い進化の歴史を通じて予測不可能な環境変化に対応できた個体のみが現在に至っていると考えられる。生物は、生きつづけるために優れた親の性質を遺伝子として子に伝える。

このような、生物進化の法則と遺伝のメカニズムを工学的に取り入れ、近年はモダンヒューリスティクス³⁾として、最適化アルゴリズムとして構成するものである⁴⁾。

3.1 遺伝アルゴリズムの概念

生物の各個体は、それぞれ固有の染色体を持ち、染色体は遺伝子の配列で構成されている。ここで、決定変数*x*を染色体に対応させて、次式のような記号列で表す。

$$x : M_1 M_2 \dots M_i \dots M_N \quad (8)$$

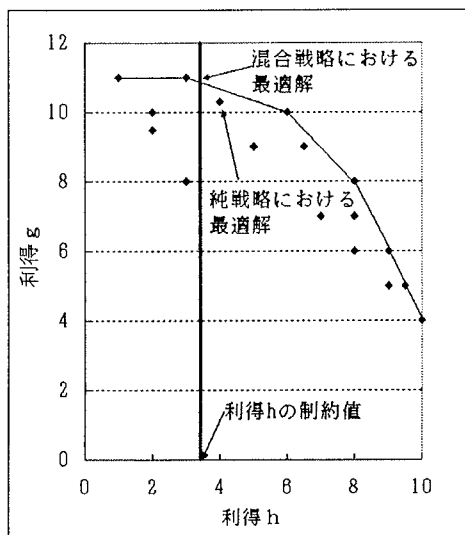


図1: 混合政策と純正策での最適化

ここに、 A_i は遺伝子に対応し、遺伝子が置かれている位置を遺伝子座と呼ぶ。また、各遺伝子が取り得る値を対立遺伝子と呼ぶ。その値は、0か1の整数、1からMまでの整数など、問題に応じて定義される。

上式のような記号列の表現を遺伝子型と呼び、その遺伝子によって定まる個体の性質を表現型と呼ぶ。

自然界における生物の進化過程では、ある世代を形成している個体の集合(個体群)を考え、この個体の中で環境への適応度の高い個体が多く生き残るように淘汰される。そして、交叉や突然変異が生じて、次の世代が構成される。これを最適化問題を解く繰返し過程に対応させる。すなわち、問題の解の候補を複数個選んでおき、第*t*回目の繰返し計算における解集合を次式のように構成する^{5) 6)}。

$$X(t) = \{x_1(t), x_2(t), \dots, x_S(t)\} \quad (9)$$

ここで、 S は個体群のサイズを表す。 S 個の解の集合である個体群は、淘汰、交叉、および突然変異という操作(遺伝演算子)を受けて、次世代の個体群を生み出す。このような操作を繰り返して、世代を十分経た後の個体群は最適解の近傍に収束すると考えられている。

step1

世代を $t=0$ とする。 S 個の個体(政策)をランダムに生成して、初期個体群 $X(0)$ 、

$$X(0) = \{x_1(0), x_2(0), \dots, x_S(0)\}$$

を設定する。(但し、各個体の遺伝子は1~ K の10進数表示。)

step2

各個体の表現型を考慮して、適応度を決める。この適応度に依存した一定のルールで個体の淘汰を行なう。(ルール戦略、エリート保存戦略、ランク戦略)

step3

一定の確率で交叉、突然変異を行い、新しい個体を生成。子は親と置き変わり新しい世代 $X(t+1)$ 、

$$X(t+1) = \{x_1(t+1), x_2(t+1), \dots, x_S(t+1)\}$$

が形成される。

step4

終了条件により終了もしくは $t=t+1$ として step2 へ戻る。

このアルゴリズムの主要部分は、適応度設定と適応度の高い個体を残す手続き、および新しい個体を生成する手続きである。すなわち、淘汰により良質な個体を重点的に固執して探索し、同時に交叉や突然変異により、解の探索空間を広げているのである。これらの手続きが有効に働く時遺伝アルゴリズムは効力を発揮するのである^{7) 8)}。

3.2 遺伝アルゴリズムの適用法

この節では、制約付きマルコフ決定過程の遺伝アルゴリズムへの導入、各パラメータの設定、及び設定した3ケースの適応度について説明する。前節における記号列で表される個体 $M_1 M_2 \dots M_1 \dots M_N$ がマルコフ決定過程における純政策にあたり、遺伝子 M_i が状態 i における決定にあたる。また個体の長さ N は状態数となり、個体の遺伝子座 i に入ることができる遺伝子の数が、状態 i で選択できる決定の数である。

以下に、本研究における遺伝アルゴリズムの適用手順について述べる。

現個体群を U とし、対象とする個体(政策)を p とする。まず、(5)、(6)より π_j^p を求め、(3)、(4)より $G(p)$ 、 $H(p)$ を計算し、表現型 (h^p, g^p) とする。主な、パラメータを以下に示す。

個体(政策) $p \in U$ の表現型: (h^p, g^p)

個体の長さ(状態数): N

個体群数: S

個体 $p(p \in U)$ の適応度: F^p

また、適応度については次の3つのケースを設定し、数値実験を行った。

<CASE 1> 政策改善法を考慮しない場合

CASE 1はGAだけの探索で、時間平均利得 h の制約値 α を満たさない個体に対しては、ペナルティーとして適応度を0にし、次世代の遺伝子として継承しないようにした。ペナルティーとして h と α の乖離度に応じて適応度を減少させることもできるが、今回はGAのみの探索で

政策4での直接利得a, 直接利得b, 推移確率(1/1000)

j	a	b	1	2	3	4	5	6	7	8	9	10
i=1	3	39	446	554	0	0	0	0	0	0	0	0
i=2	8	38	436	221	344	0	0	0	0	0	0	0
i=3	19	48	0	620	315	65	0	0	0	0	0	0
i=4	0	5	0	0	196	297	507	0	0	0	0	0
i=5	10	48	0	0	0	127	253	620	0	0	0	0
i=6	5	25	0	0	0	0	240	625	135	0	0	0
i=7	11	46	0	0	0	0	0	347	189	463	0	0
i=8	3	35	0	0	0	0	0	0	301	610	89	0
i=9	4	44	0	0	0	0	0	0	0	109	341	551
i=10	15	35	0	0	0	0	0	0	0	0	673	327

政策5での直接利得a, 直接利得b, 推移確率(1/1000)

j	a	b	1	2	3	4	5	6	7	8	9	10
i=1	6	28	549	451	0	0	0	0	0	0	0	0
i=2	6	15	30	626	343	0	0	0	0	0	0	0
i=3	6	40	0	559	235	206	0	0	0	0	0	0
i=4	13	14	0	0	264	595	142	0	0	0	0	0
i=5	16	12	0	0	0	258	379	363	0	0	0	0
i=6	11	26	0	0	0	0	409	181	409	0	0	0
i=7	18	26	0	0	0	0	0	316	377	307	0	0
i=8	11	11	0	0	0	0	0	0	156	269	575	0
i=9	18	15	0	0	0	0	0	0	0	215	355	430
i=10	15	29	0	0	0	0	0	0	0	0	603	397

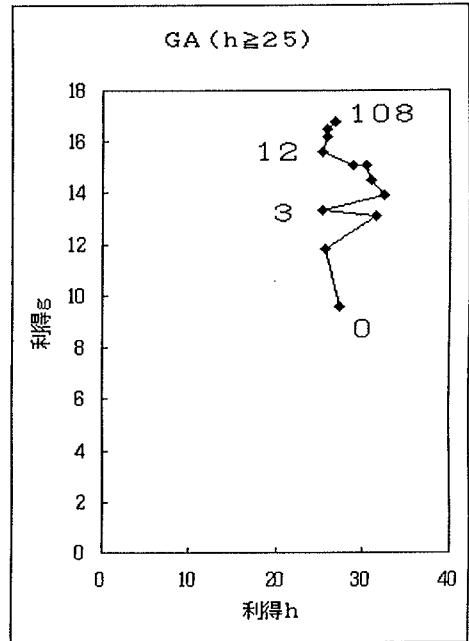


図3: CASE1($h > 25$)での(h,g)の変化

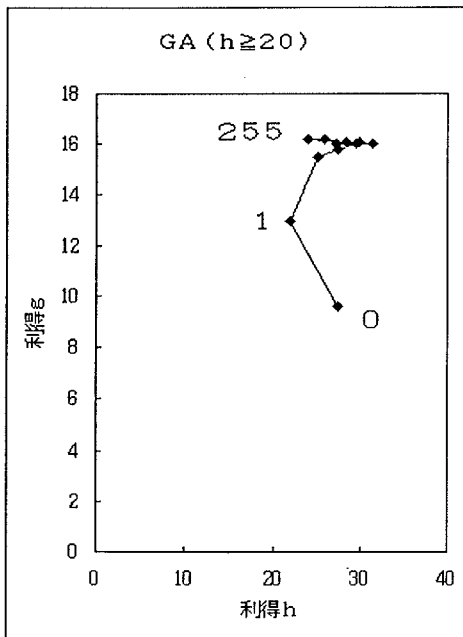


図2: CASE1($h > 20$)での(h,g)の変化

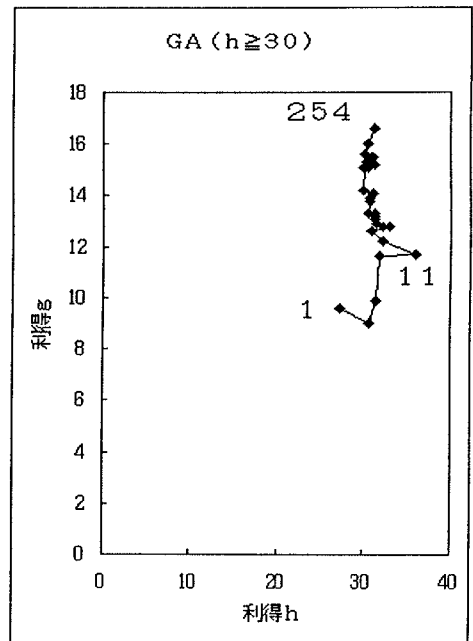


図4: CASE1($h > 30$)での(h,g)の変化

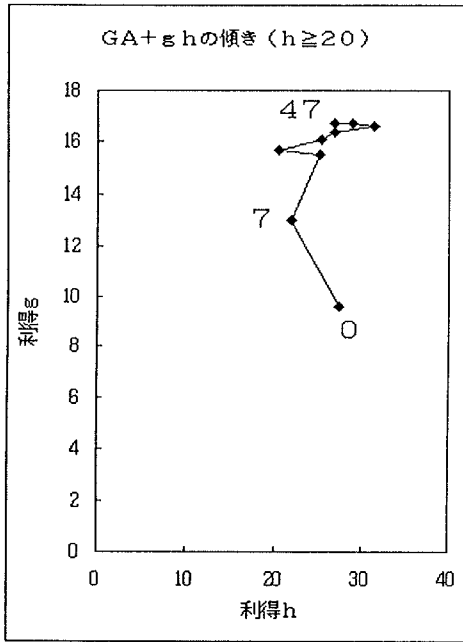


図 5: CASE2($h > 20$) での (h, g) の変化

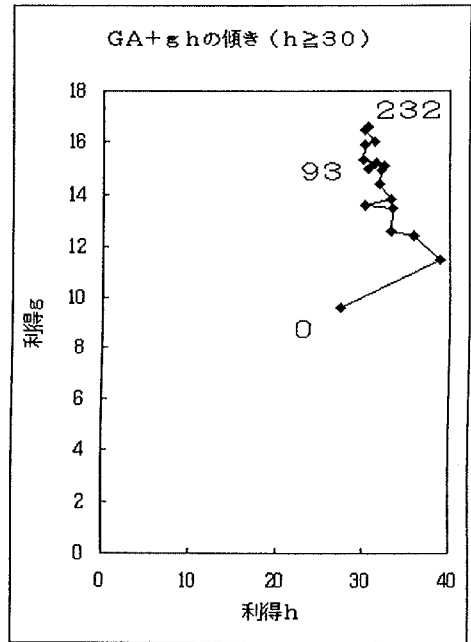


図 7: CASE2($h > 30$) での (h, g) の変化

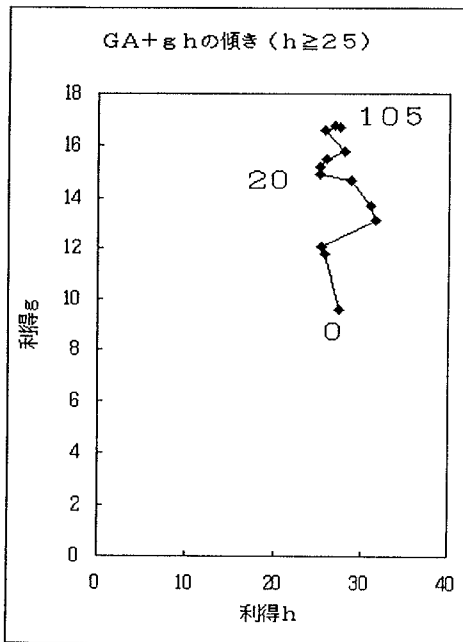


図 6: CASE2($h > 25$) での (h, g) の変化

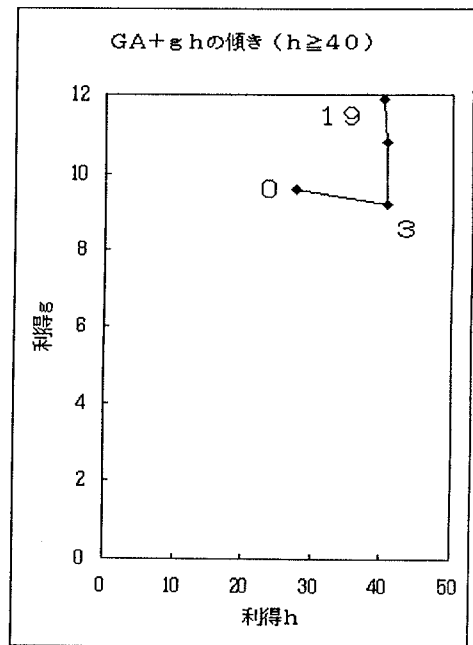


図 8: CASE2($h > 40$) での (h, g) の変化

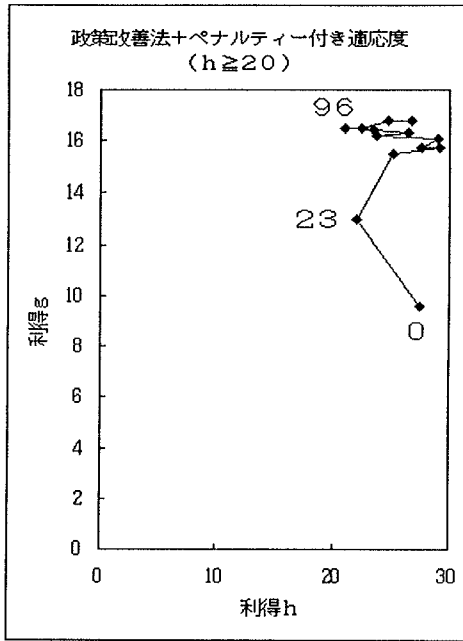


図 9: CASE3($h > 20$)での (h,g) の変化

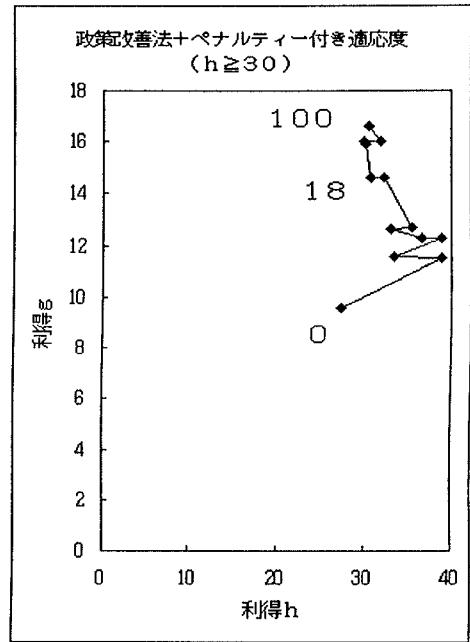


図 11: CASE3($h > 30$)での (h,g) の変化

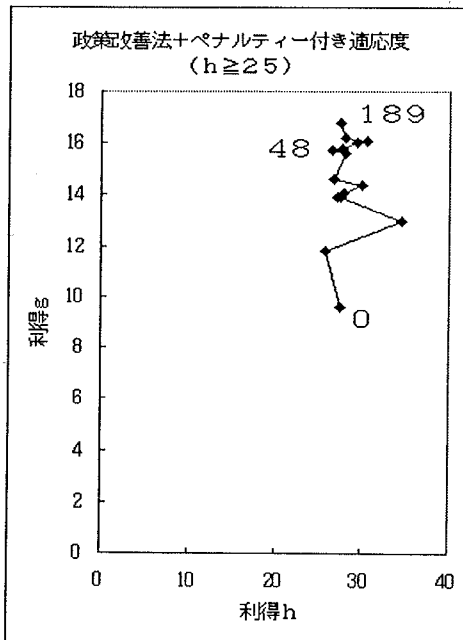


図 10: CASE3($h > 25$)での (h,g) の変化

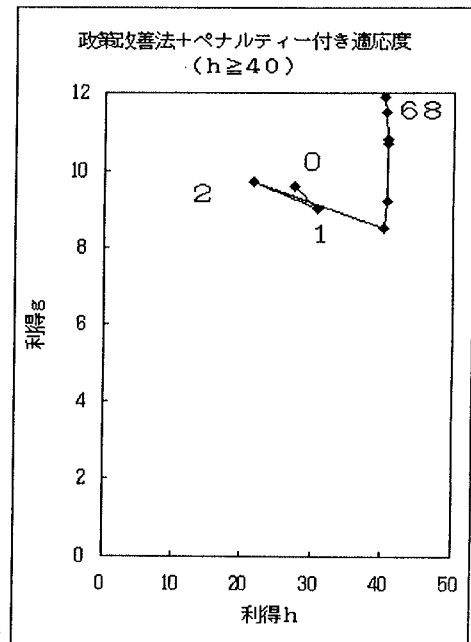


図 12: CASE3($h > 40$)での (h,g) の変化

前節で提案した3つのCASEについて、時間平均利得 h の制約値 α を変化させて、数値計算を行ったのでその結果を示す。これらの数値計算は全て同じ初期解でいずれも300世代まで計算した結果である。図2~図4は α が20、25、30、40のときのCASE1での世代推移における (h, g) の値の変化を示したものである。図中の数字は世代数を示している。図5~図8はCASE2での世代推移における (h, g) の値の変化を示したものである。図9~図12はCASE3での世代推移における (h, g) の値の変化を示したものである。

5 考察

CASE1のGAだけの探索では予想以上の効果があった。しかし、制約値 α の値が大きくなるにつれ、制約を満たした解を見つけるまでに時間がかかっている。また $\alpha=40$ の時には300世代でも制約を満たした解を探索することができなかった。

CASE2のGAと政策改善法のハイブリッド型では、CASE1よりも早期に制約を満たした解を探索していることが判る。また、CASE1では探索不可能であった $\alpha=40$ の時でもわずか18世代で最適解に到達している。

CASE3のハイブリッド型+適応度ペナルティーを与えるGAでは、CASE1の約半分の世代でCASE1同等もしくはそれ以上の探索能力を発揮している。また、 $\alpha=20$ の時にはCASE2の方が早く最適解に到達しているように見えるが実はCASE2は最適解には到達してはならず g の値は16.7であった。しかし、CASE3では最適解 $g=16.8$ ($\alpha=25$ の時と同じ)に達していた。また、300世代以内でCASE3は α の値がいずれの時も最適解に達していた。

これらのことから、GAだけのランダム探索よりもGAと政策改善法のハイブリッド型で構成した探索法の方が効率的な探索が実現できていることが判る。

時間平均利得 h の制約値 α が大きくなるほど、探索空間は小さくなり、実行可能解でさえ探索は困難になる。逆に、制約値 α が小さくなるほど、探索空間は広がり、実行可能解の中から最適解を探索することが困難となる。今回の数値実験ではどちらの場合でも政策改善法とGAのハイブリッド型がGAだけの探索よりも有効であることが確認された。

また、前者の場合にはCASE2の適応度を政策改善法によって更新される G の増分とした方が有効であり後者ではCASE3の適応度にペナルティーを与える方が有効であろう。これは、制約を満たしていない個体(政策)の適応度にペナルティーを与えることにより、個体群内に無駄な探索となる個体を留めないためであると考えられる。

6 おわりに

本研究では制約付きMDP問題について、GAと政策改善法のハイブリッド型を提案したが、非常に良い結果が得られた。今回は政策改善法は各個体(政策)に1回しか行っていないが、繰り返し行えば必ず制約を満たす個体を生成することも可能である。これは次回の課題としたい。

適応度の設定方法の違いによって、同じハイブリッド型でもCASE2、CASE3のように探索過程が異なってくることは興味深い。また、適応度の設定方法は今回の数値実験を行った方法以外にも、様々な方法が考えられる。

制約についても、今回は1つであったがGAでは複数の制約も取扱うことが可能である。しかし、その際には適応度の設定方法をよく考慮しておかないと効率的な探索は行えないであろう。

今後、これらの課題についても研究を継続していきたい。

参考文献

- 1) Beutler, F.J. and Ross, K.W.: Optimal Policies for Controlled Markov Chains with a Constraint, *J. Math. Anal. Appl.*, Vol. 112, pp. 236-252, 1985.
- 2) H. Kawai, N. Katoh: Variance Constrained Markov Decision Process, *Journal of Operations Research Society of Japan*, Vol. 30 No. 1 March 1987
- 3) 北川敏夫: マルコフ過程、共立出版
- 4) 茨木俊秀: 組合せ最適化法をめぐる最近の話題、モダンヒューリスティックスの新展開 - Genetic Algorithm, Simulated Annealing, Tabu Search, Neural Net 法は本当に有効か? -、日本オペレーションズ・リサーチ学会第30回シンポジウム, pp. 1-10(1993).
- 5) 北野宏明: 遺伝アルゴリズム、産業図書、(1993)
- 6) 樋口哲也、北野宏明: 遺伝アルゴリズムとその応用、情報処理 July 1993 Vol. 34 No. 7 p. 871~p. 883
- 7) 三宮信夫: 遺伝アルゴリズムによる最適化問題の解法、第36回システム制御情報学会研究発表公演会 p. 9~p. 18
- 8) Branko, Soucek, and The IRIS Group: DYNAMIC, GENETIC, AND CHAOTIC PROGRAMING, WILEY INTER SCIENCE.
- 9) David. E. Goldberg: Genetic Algorithms in Search Optimization and Machine Learning, Wesley Publishing Company INC(1989).