

## ニューロ DP による多品目在庫管理の最適化

愛知工業大学 大野勝久(Katsuhisa Ohno)  
Aichi Institute of Technology,  
名古屋工業大学 伊藤崇博(Takahiro Ito)  
Nagoya Institute of Technology,  
大阪府立大学 石垣智徳(Tomonori Ishigaki)  
Osaka Prefecture University,  
NTT ドコモ東海 渡辺 誠(Makoto Watanabe)  
NTT DoCoMo Tokai

## 1. 序論

在庫管理問題は、古くは 1913 年の F.W.Harris の経済的発注量(EOQ)の研究に始まり、経営工学分野の中で長年にわたり中心的課題として、多くの研究が行われてきた。今日においては、原材料の調達から消費者にいたる SCM を中心に、多品目在庫管理問題の最適化が切望されている。

本発表ではまず、2 章で離散時間多品目在庫管理問題を平均費用を最小化する時間平均マルコフ決定過程問題として定式化する。動的計画法(以下 DP と略す)の枠組みとマルコフ決定過程(MDP)は、Bellman によって 1950 年代に提案されたが、政策反復法(Policy Iteration Method, PIM)を初めとするアルゴリズムは Howard[1], Puterman[2]等にまとめられている。離散時間多品目在庫管理問題にたいしても、Ohno et al.[3]は、PIM を改善したアルゴリズムを提案している。しかし、MDP は DP 同様、品目数の増加とともに次元の呪いを引き起こし、実用規模の問題を事実上解くことができない。そこで本研究では、人工知能[4]の分野において強化学習(Reinforcement Learning)[5]とも呼ばれている、ニューロ・ダイナミックプログラミング(Neuro-Dynamic Programming, ニューロ DP) [6, 7]の適用を試みる。ニューロ DP は MDP の枠組みの中で、シミュレーション、学習、ニューラルネットワークなどを組み合わせ、大規模な問題に対する近似最適政策を計算する手法として開発されてきた。最近 Das et al.[8]と Gosavi[9]は、MDP を含むより一般的なセミ・マルコフ決定過程(SMDP)の時間平均費用を最小化するアルゴリズムとして SMART, RELAXED-SMART を提案し、保全問題へ適用している。また He et al.[10]は、PIM の値決定ルーチンをシミュレーションで置きかえた SBPI アルゴリズムを提案し、在庫管理問題へ適用している。さらに Gosavi et al.[11]は、SMART に TD(Temporal Difference) をとり入れた  $\lambda$ -SMART を開発し、

航空運賃の収入管理問題へ応用している。

本発表では、まず Ohno et al.[3]に基づき次章で多品目在庫管理問題を説明し、3 章で時間平均マルコフ決定過程(UMDP)への定式化を行い、若干の理論的性質を定理として示す。4 章で修正政策反復法(Modified Policy Iteration Method, MPIM) [12-14]の値決定ルーチンにシミュレーションを用いる SBMPIM アルゴリズム[15]を述べる。5 章では SMART, RELAXED-SMART と SBPI アルゴリズムを簡単に説明する。6 章ではこれらアルゴリズムを 2 品目在庫管理問題へ適用し、SBMPIM, SMART, RELAXED-SMART, SBPI に対する結果をも含めて Ohno et al.[3]による厳密解との比較を行う。

## 2. 多品目在庫管理問題

発注費用の形が一般的であり、超過需要が部分的もしくはすべてバックログされ、周期的な観測を行う  $N$  品目在庫管理システムを考える。もしロスセールを仮定するならば、品物の納入遅れはないものとする。ここで  $\mathbf{x} = (x_1, x_2, \dots, x_N)^T$  を各期首における発注前の手持ち在庫水準とする。ここで  $x_n \in \mathbf{Z} = \{0, \pm 1, \pm 2, \dots\}$  ( $n \in \mathbf{N} = \{1, 2, \dots, N\}$ ) であり、 $T$  は転置を表す。Johnson[16]を除くすべての論文は連続状態を論じているが、より現実的な離散状態を議論する。発注は  $\mathbf{x}$  と発注政策  $f$  によって決定され、 $\mathbf{y} = (y_1, y_2, \dots, y_N)^T \in \mathbf{Z}^N$  まで発注したとき、段取り費用を含む発注費  $K(\mathbf{x}, \mathbf{y})$  が発生する。各期の需要  $\mathbf{D} = (D_1, D_2, \dots, D_N)^T$ ,  $D_n \in \mathbf{Z}_+ = \{0, 1, 2, \dots\}$  ( $n \in \mathbf{N}$ ) は独立な同一分布  $\phi(\mathbf{y}, \mathbf{D})$  に従うものと仮定する。その期に需要  $\mathbf{D}$  が発生したならば、1 期間当たり在庫・品切れ費用  $I(\mathbf{y}, \mathbf{D})$  が発生する。したがって、各期における期待在庫・品切れ費用  $L(\mathbf{y})$  は

$$L(y) = \sum_{D \geq 0} \phi(y, D) l(y, D)$$

で与えられる。ここで  $x \leq y$  ( $x, y \in Z^N$ ) は  $x_n \leq y_n$  ( $n \in N$ ) を表し、 $x < y$  は  $x \leq y$  かつ  $x \neq y$  を意味するものとする。また、次期の期首在庫は  $u(y, D)$  によって与えられるものとする。例えば、バックログを仮定する場合、 $u_B(y, D) = y - D$  となり、ロスセールを仮定する場合、

$$u_L(y, D) = ([y_1 - D_1]^+, [y_2 - D_2]^+, \dots, [y_N - D_N]^+)^T, \\ [y_n - D_n]^+ = \max\{y_n - D_n, 0\} \quad (n \in N)$$

である。

状態  $x$  の値からなる状態空間  $X$  は倉庫容量  $\bar{M}$  と在庫処理環境 (例えばバックログ (B), ロストセール (L)) によって制約される。例えば、バックログの場合

$$X_B = \left\{ x \in Z^N; \sum_{n \in N} a_n x_n^+ \leq \bar{M} \right\}$$

であり、ロスセールの場合

$$X_L = \left\{ x \in Z_+^N; \sum_{n \in N} a_n x_n \leq \bar{M} \right\}$$

となる。ここで  $a_n$  は  $n$  番目の品目 1 個当たりの大きさであり、 $x_n^+ = \max\{x_n, 0\}$ ,  $n \in N$  である。 $X_B$  は加算集合であり  $X_L$  は有限集合である。以下の条件を仮定する。

条件 1  $x < y < z$  であるような任意の在庫水準  $x, y, z \in X$  に対して以下の式が成立する。

$$K(x, z) < K(x, y) + K(y, z) \text{ かつ } K(x, x) = 0.$$

条件 2 もし  $X$  が有限でないならば  $\lim_{x \rightarrow \infty} L(x) = \infty$  である。すなわち、総ての実数  $M$  に対して、 $x \in X$  かつ  $x \geq z_M$  を満たすすべての  $x$  に対して  $L(x) > M$  となる  $z_M \in X$  が存在する。

条件 3 もし  $X$  が有限でないならば各期の需要は有界である。すなわち、すべての  $x \in X$  に対して

$$\sum_{D \in B} \phi(x, D) = 1$$

を満たすような有界集合  $B \subset Z_+^N$  が存在する。

条件 4 ロストセールの仮定の下での  $x_n = 0$  を除くすべての  $x \in X$  と  $n \in N$  に対し  $x_n > u_n(x, D)$  が確率 1 で成立する。

Johnson[16]では固定発注費用  $K > 0$  に対して

$$K(x, y) = K\delta(y - x) + c \cdot (y - x)$$

と仮定されている。ここで  $\delta(0) = 0$  であり、 $x > 0$  ならば  $\delta(x) = 1$ ,  $c$  は  $N$  次元の価格行ベクトルである。明らかにこの発注費用は条件 1 を満たし、条件 2 と条件 3 は可算である  $X$  に対してのみ要

求される条件である。また、確率 1 で需要 0 をとる品目は在庫管理の対象として無視できるので条件 4 は一般的な在庫管理システムについて成立する。目的は計画期間が無限にわたる在庫管理システムの単位期間当たりの期待費用を最小にする最適発注政策を見つけることである。

### 3. マルコフ決定過程による定式化

前章の多品目在庫管理システムは UMDP に定式化することができ、その最適方程式は次式で与えられる。

$$g(x) + v(x) = \min_{y \in Y(x)} \left\{ K(x, y) + L(y) + \sum_{D \in B} \phi(y, D) v(u(y, D)) \right\} \quad (1)$$

$$g(x) = \min_{y \in Y} \left\{ \sum_{y \in Y} \phi(y, D) g(u(y, D)) \right\} \quad (2)$$

ここで

$$Y = \{y \in X; y \geq x\},$$

$$Y(x) = \left\{ y \in Y; g(x) = \sum_{D \in B} \phi(y, D) g(u(y, D)) \right\}.$$

である。 $g(x)$ ,  $v(x)$  をそれぞれ初期状態が  $x$  である時の最小期待平均費用と相対値とおく。有限である  $X$  の場合、確定的な定常マルコフ政策の中に最適政策が存在することがよく知られている。定常政策  $f$  は  $X$  の部分集合  $\sigma$  と  $\sigma$  から  $X$  への写像  $S(\cdot)$  を用いて  $f = (\sigma, S(\cdot))$  と表わすことができる。 $\sigma$  は  $(\sigma, S(\cdot))$  政策は在庫水準  $x$  において次のように発注を決定する。

$x \in \sigma$  のとき  $S(x)$  まで発注する。

$x \in \sigma^c \equiv X - \sigma$  のとき 発注しない。

したがって、 $\sigma$  を発注状態集合、 $\sigma^c$  を非発注状態集合と呼ぶことにする。

定理 1 最適発注政策  $(\sigma, S(\cdot))$  は以下を満たす。

i) もし  $X$  が有限でないならば

$$\sigma^c \subset (\sigma^0)^c = \{x \in X; x \geq z_g^0\}, \quad (3)$$

である。ここで  $z_g^0$  は仮定 2 の  $M$  を以下で与えられる  $g^0$  によって置き換えたものである。

$$g^0 = \min_{y \in X} \left\{ L(y) + \sum_{D \in B} \phi(y, D) K(u(y, D), y) \right\}. \quad (4)$$

ii)  $S(x) \in \sigma^c$ ,  $x \in \sigma$ .

iii)  $v(x) = K(x, S(x)) + v(S(x))$ ,  $x \in \sigma$ .

(5)

もし  $\mathbf{X}$  が有限でないならば  $\mathbf{X}^0$  を

$$(\sigma^0)^c \cup \left\{ u(\mathbf{x}, \mathbf{D}) \text{ for all } \mathbf{x} \in (\sigma^0)^c \text{ かつ } \mathbf{D} \in \mathbf{B} \right\}$$

とおき、 $\mathbf{X}^0$  の状態数を  $N^0$  とする。定理 1 は最適政策の候補となるようなすべての発注政策の下で  $\mathbf{X}^0$  がすべての再帰的状态を含むことを意味している。逆に  $\mathbf{X}^0$  の外側のすべての状態は最適政策の下で常に過渡的状态であり、この状態における発注政策は最小平均期待費用に影響しない。したがって、有限集合  $\mathbf{X}^0$  のみに対して最適政策を探せば十分である。 $\mathbf{x} \notin \mathbf{X}^0$  に対する最適政策  $S(\mathbf{x})$  は(5)式より容易に計算できる。次の定理 2 では  $\mathbf{X}$  が有限であるとき、 $\mathbf{X}^0$  は  $\mathbf{X}$  を意味するものとする。

定理 2 最小平均期待費用  $g(\mathbf{x})$  は  $\mathbf{x} \in \mathbf{X}$  に対して定数である。すなわち、 $g(\mathbf{x}) = g$  であり、最適方程式(1)と(2)は次式に帰着する。

$$g + v(\mathbf{x}) = \min_{\mathbf{y} \in \mathbf{Y}} \left\{ K(\mathbf{x}, \mathbf{y}) + L(\mathbf{y}) + \sum_{\mathbf{D} \in \mathbf{B}} \phi(\mathbf{x}, \mathbf{D}) v(u(\mathbf{x}, \mathbf{D})) \right\},$$

$$\mathbf{x} \in \mathbf{X}^0, \quad (6)$$

ここで右辺の最小値を与える決定が最適政策である。

#### 4. 修正政策反復法とニューロ DP アルゴリズム

最適性方程式(6)を解くアルゴリズムが PIM, MPIM であり、値反復法 (value iteration method), 線形計画法 (linear programming) である[1,2,12-14]。以下、表記を簡単化するため

$$r(\mathbf{x}, \mathbf{y}) = K(\mathbf{x}, \mathbf{y}) + L(\mathbf{y})$$

$$\mathbf{x}' = u(\mathbf{y}, \mathbf{D})$$

とおく。

[MPIM] [14]

ステップ 1:  $v_0(s_r) = 0$  をみたく初期ベクトル  $v^0$ , 非負整数  $m$ , 初期政策  $f^0$ , 正数  $\varepsilon$  を定め、 $k=0$  とおく。

ステップ 2: (政策改良ルーチン) 各  $\mathbf{x} \in \mathbf{X}^0$  に対して

$$g^{k+1}(\mathbf{x}) = \min_{\mathbf{a} \in K(\mathbf{x})} \left\{ r(\mathbf{x}, \mathbf{y}) + \sum_{\mathbf{D} \in \mathbf{B}} \phi(\mathbf{x}, \mathbf{D}) v^k(\mathbf{x}') - v^k(\mathbf{x}) \right\}$$

を計算し、 $f^k(\mathbf{x})$  が  $g^{k+1}(\mathbf{x})$  を与えれば、 $f^{k+1}(\mathbf{x}) = f^k(\mathbf{x})$  とおき、さもなければ、 $g^{k+1}(\mathbf{x})$  を与える任意の決定を  $f^{k+1}(\mathbf{x})$  ととる。

ステップ 3: (値近似ルーチン)

$$w^0(\mathbf{x}) = v^k(\mathbf{x}) + v^{k+1}(\mathbf{x}), \quad \mathbf{x} \in \mathbf{X}^0$$

とおき、 $l=0, 1, \dots, m-1$  に対して順次、

$$w^{l+1}(\mathbf{x}) = r(\mathbf{x}, f^{k+1}(\mathbf{x})) + \sum_{\mathbf{D} \in \mathbf{B}} \phi(\mathbf{x}, \mathbf{D}) w^l(\mathbf{x}'),$$

$\mathbf{x} \in \mathbf{X}^0$

を計算し、

$$v^{k+1}(\mathbf{x}) = w^m(\mathbf{x}) - w^m(\mathbf{x}_r), \quad \mathbf{x} \in \mathbf{X}^0$$

とおく。すべての  $\mathbf{x}$  に対して、 $|v^{k+1}(\mathbf{x}) - v^k(\mathbf{x})| < \varepsilon$

であれば終了。最適政策は  $f^{k+1}(\mathbf{x})$ , 最小平均費用は  $g^{k+1}(\mathbf{x}_r)$ , 相対費用は  $v^{k+1}(\mathbf{x})$  で与えられる。さもなければ、 $k=k+1$  として、ステップ 2 へ。

MPIM は比較的状态数の大きな MDP に対しても有効であるが、DP 同様、品目数の増加とともに次元の呪いを引き起こし、多品目問題を解くことは困難である。したがって、MPIM において状態空間  $\mathbf{X}^0$  の全ての状態にたいして値近似ルーチンを実行することは実際的ではなく、シミュレーションを用いることが考えられる。すなわち、実際によく生起する初期状態  $\mathbf{x}_0$  から出発し、システムの状態変化と費用をシミュレートし、訪問した状態の集合  $\mathbf{X}_v$  にたいしてだけ相対費用  $v(\mathbf{x})$  を推定する。このニューロ DP アルゴリズムを SBMPIM (Simulation-Based Modified Policy Iteration Method) と呼ぶことにする[15]。

[SBMPIM] [15]

ステップ 1: 初期状態  $\mathbf{x}_0$  と望ましい状態  $\mathbf{x}^*$  を定め、シミュレーション回数  $m$  および  $\lambda$  ( $0 \leq \lambda \leq 1$ ) を定めて、訪問した状態の集合  $\mathbf{X}_v = \mathbf{X}_T = \phi$  (空集合), 累積費用  $TC = 0$ ,  $\mathbf{x} = \mathbf{x}_0$ ,  $k=l=1$  とおく。

ステップ 2:  $\mathbf{x} \notin \mathbf{X}_v$  ならば、 $\mathbf{X}_v = \mathbf{X}_v \cup \{\mathbf{x}\}$ ,  $\mathbf{X}_T = \mathbf{X}_T \cup \{\mathbf{x}\}$ ,  $\mathbf{x}$  の訪問回数  $v(\mathbf{x}) = 1$  とおき、 $f(\mathbf{x})$  を状態  $\mathbf{x}^*$  へ向かう実行可能な決定と定める。 $\mathbf{x} \in \mathbf{X}_v$  ならば、 $\mathbf{x} \notin \mathbf{X}_T$  のとき、 $\mathbf{X}_T = \mathbf{X}_T \cup \{\mathbf{x}\}$ ,  $v(\mathbf{x}) = 1$  とおき、 $\mathbf{x} \in \mathbf{X}_T$  ならば、

$$v(\mathbf{x}) = v(\mathbf{x}) + 1$$

と更新する。状態  $\mathbf{x}$  で決定  $f(\mathbf{x})$  をとったときの状態推移をシミュレーションし、次期の状態  $\mathbf{x}'$  を定める。

$$TC = TC + r(\mathbf{x}, f(\mathbf{x}))$$

$$\mathbf{x} = \mathbf{x}'$$

と更新し、 $l=m$  ならばステップ 3 へ。さもなければ  $l=l+1$  としてステップ 2 へ。

ステップ 3: ( $g$  の推定) 平均費用  $g$  を次式により推定する。

$$g = TC/m$$

ステップ 4: ( $v(\mathbf{x})$  の推定)  $\mathbf{X}_v$  のなかで  $\mathbf{x}_r$  を定め、 $v(\mathbf{x}_r) = (1 - \lambda v(\mathbf{x}_r))/m (w(\mathbf{x}_r) - g) + (\lambda v(\mathbf{x}_r))/m (r(\mathbf{x}_r, f(\mathbf{x}_r)) - g)$  を計算し、 $\mathbf{x} (\neq \mathbf{x}_r) \in \mathbf{X}_v$  にたいして

$v(x) = (1 - \lambda v(x)/m)(w(x) - g) + (\lambda v(x)/m)(r(x, f(x)) - g) - v(x)$  を計算し,  $v(x_r) = 0$  とおく。ただし,  $k=1$  のときには  $v(x_r) = r(x_r, f(x_r)) - g$ ,  $v(x) = r(x, f(x)) - g - v(x_r)$  である。

ステップ 5: (政策改良ルーチン)  $x \in X_v$  にたいして

$$w(x) = \min_{y \in N(x, f(x))} \left\{ r(x, y) + \sum_{D \in B} \phi(y, D) v(x') \right\}$$

を計算し,  $v(x) = 0$  とおく。ここで  $N(x, f(x))$  は  $f(x)$  の近傍であり,  $\phi(y, D) > 0$  となる  $x \in X_v$  にたいしては,  $X_v = X_v \cup \{x\}$ ,  $v(x') = 0$  とおき,  $f(x')$  を  $x^*$  へ向かう実行可能な決定と定める。  
 $w(x') = r(x', f(x'))$  とおき,

$$v(x') = r(x', f(x')) - r(x_r, f(x_r))$$

として,  $w(x)$  を計算する。 $f(x)$  が  $w(x)$  を与えなければ,  $w(x)$  を与える任意の決定として  $f(x)$  を改良する。 $k$  が停止回数に達すれば終了。さもなければ  $X_T = \phi$ ,  $TC = 0$ ,  $l = 1$ ,  $k = k + 1$  とおきステップ 2 へ。

## 5. SMART と SBPI アルゴリズム

MPIM と SBMPIM を最適制御問題に適用するに先立ち, 既存の NDP アルゴリズムを簡単に紹介する。

[SMART] [8]

ステップ 1: 全ての  $x \in X$  と  $y \in K(x)$  にたいして Q-factor  $Q_{new}(x, f) = Q_{old}(x, f) = 0$ , 累積費用  $TC = 0$ , 累積時間  $T = 0$ , 平均費用  $g = 0$ , 反復回数  $k = 0$  とおき, 学習率 (learning rate) と探検確率 (exploration probability) をステップ 2 で定める定数  $(\alpha_0, \alpha_\tau, p_0, p_\tau)$  を与える。

ステップ 2: 反復  $k$  で状態  $x$  にいれば, 学習率  $\alpha_k$ , 探検確率  $p_k$  を

$$\alpha_k = \alpha_0(\alpha_\tau + k) / (k^2 + k + \alpha_\tau),$$

$$p_k = p_0(p_\tau + k) / (k^2 + k + p_\tau)$$

として定める。

ステップ 3: 高い確率  $(1 - p_k)$  で  $Q_{new}(x, y)$  を最小にする決定  $y^*$  を選択し, 確率  $p_k$  で  $y^*$  を除く  $K(x)$  からランダムに  $y$  を選択する。

ステップ 4: 選択された決定  $y$  でシミュレーションを行い, 状態  $x'$  へ推移すれば, 直接費用  $r(x, x', y)$  がかかる。

ステップ 5:  $Q_{new}(x, y)$  を次式により更新する。

$$Q_{new}(x, y) = (1 - \alpha_k) Q_{old}(x, y) + \alpha_k \left\{ r(x, x', y) - g + \min_{y' \in K(x')} Q_{old}(x', y') \right\}$$

ステップ 6: ステップ 3 で決定  $y^*$  を選択したならば,  $TC$  と  $g$  を更新する。

$$TC = TC + r(x, x', y^*)$$

$$T = T + 1$$

$$g = TC/T$$

ステップ 7:  $Q_{old}(x, y) = Q_{new}(x, y)$  と更新する。

ステップ 8:  $k$  が停止回数に達すれば終了。さもなければ  $k = k + 1$ ,  $x'$  を  $x$  としてステップ 2 へ。

Gosavi[9]は SMART が必ずしも収束しないことを示し, その改良版 RELAXED-SMART を提案している。SMART と RELAXED-SMART は SMDP にたいするアルゴリズムであるが, ここでは MDP にたいするものに修正している。

[RELAXED-SMART] [9]

ステップ 3~5, 7, 8 は [SMART] と同じである。

ステップ 1: Q-factor  $Q_{new}(x, y) = Q_{old}(x, y) = 0$ ,  $TC = 0$ ,  $T = 0$ ,  $g = 0$ ,  $k = 0$  とおき, パラメータ  $\alpha_0$ ,  $p_0$ ,  $\beta_0$  を与える。

ステップ 2: 反復  $k$  で状態  $x$  にいれば,  $\alpha_k$ ,  $p_k$ ,  $\beta_k$  を

$$\alpha_k = \alpha_0/k, \quad p_k = p_0/m, \quad \beta_k = \beta_0/k$$

として定める。

ステップ 6: ステップ 3 で決定  $y^*$  を選択したならば,  $TC$ ,  $T$ ,  $g$  を次式で更新する。

$$TC = (1 - \beta_k)TC + \beta_k r(x, x', y^*)$$

$$T = (1 - \beta_k)T + \beta_k$$

$$g = TC/T$$

一方, He et al.[10]は PIM の値決定ルーチンをシミュレーションで置きかえた SBPI (Simulation Based Policy Iteration) アルゴリズムを提案している。

[SBPI アルゴリズム] [10]

ステップ 1: 初期政策  $\{f_0(x); x \in X\}$  を定め,  $k = 0$  とおく。

ステップ 2: (値決定ルーチン)

2-a: ( $g^k$  の推定)

i) 初期状態  $x_0$  からシミュレーションにより  $x_1, \dots, x_m$  を生成する。

ii)  $g^k = 0$  とおき,  $n = 0, \dots, m-1$  にたいして  $(x_n, x_{n+1})$  の推移に伴う  $g^k$  を次式で更新する。

$$g^k = (1 - 1/(n+1))g^k + (1/(n+1))r(x_n, x_{n+1}, f^k(x_n))$$

2-b: ( $v^k(x)$  の推定)

i) ステップ 2-a で行ったシミュレーションにおいて訪問回数が最も多い状態を  $x^*$  とおく。

ii) 過渡状態  $x_0$  から出発し, 状態  $x^*$  へ至るトラジェクトリーをシミュレーションにより  $L$  本生成する。

iii) 1 本目のトラジェクトリ  $(x_0, x_1, \dots, x_n = x^*)$ ,  $l=1, \dots, L$ , にたいして, 推移  $(x_n, x_{n+1})$  に伴う  $w(x_i)$ ,  $i=1, \dots, n$ , を次式により更新する。

$$w(x_i) = w(x_i) + \gamma_i \lambda^{n-i} d_n$$

ここで,  $\gamma_i$  はそのトラジェクトリ中で  $s_i$  を訪問した回数の逆数であり,  $0 \leq \lambda \leq 1$ ,  $d_n = r(x_n, x_{n+1}, f^k(x_n)) - g^k + w(x_{n+1}) - w(x_n)$  である。

iv)  $v^k(x) = w(x) - w(x_r)$ ,  $x \in X$

ステップ 3: (政策改良ルーチン)

$$f^{k+1}(x) = \arg \min_{y \in K(x)} \left\{ r(x, y) + \sum_{D \in B} \phi(y, D) v^k(x') \right\},$$

$x \in X$

ステップ 4:  $f^{k+1}(x) = f^k(x)$ ,  $x \in X$  ならば停止。

最適政策は  $f^k(x)$  である。さもなければ  $k = k+1$  としてステップ 2 へ。

## 6. 数値実験

各アルゴリズムは VB (Visual Basic) によりプログラムし, DOS/V 機 (CPU: Pentium4 2.4GHz, メモリ: 2GB) を用いて計算を行った。

### 6.1 数値例

- ・各品目の容積はすべて 1
- ・品目数 ; 2
- ・在庫容量 ; 2 3
- ・一期あたりの  $i$  番目の需要  $D_i$  ( $i=1, 2$ ) の確率  $\phi_i(D_i)$  ;

$$\phi_1(0) = 1/27, \phi_1(1) = 2/9, \phi_1(2) = 4/9, \\ \phi_1(3) = 8/27$$

$$\phi_2(0) = 1/16, \phi_2(1) = 1/4, \phi_2(2) = 3/8, \\ \phi_2(3) = 1/4, \phi_2(4) = 1/16$$

- ・発注費用 ;

$$K(x, y) = 15\delta(y-x) + 10\delta(y_1-x_1) + 5\delta(y_2-x_2) \\ + 6(y_1-x_1) + 3(y_2-x_2)$$

- ・保管・品切れ費用 ;

$$l(y, D) = 2[y_1 - D_1]^+ + 2[y_2 - D_2]^+ + 2[D_1 - y_1]^+ \\ + 14[D_2 - y_2]^+$$

### 6.2 SBMPIM と従来のアルゴリズムとの比較

SBMPIM および既存のアルゴリズムのパラメータを表 1 に示す。表 1 のパラメータのもとで各アルゴリズムにより計算したときの  $g$  および計

算時間を表 2 に示す。

表 3 は MPIM による最適政策である。また, 訪問回数の多い状態における MPIM および各アルゴリズムの政策との比較を表 4 に示す。表中「一致」の列には MPIM と政策が等しい場合に  $\circ$ , そうでない場合に  $\times$  を記入している。また, 「政策一致数」の行は全 300 状態中で MPIM と政策が一致した状態数を示している。

表 1 アルゴリズムに使用するパラメータ一覧

アルゴリズム	パラメータ
SBMPIM	$m=10000, \lambda=0.1,$ $x_0=(11,12), x^*=(8,9),$ 停止回数=50
SMART	$m=1000000, \alpha_0=0.9, p_0=0.3,$ $\alpha_r=10000, p_r=10000$
RELAXED-SMART	$m=1000000, \alpha_0=0.9, p_0=0.3$
SBPI	$m=10000, L=1, \lambda=0.9$

表 2 各アルゴリズムの  $g$  および計算時間

アルゴリズム	$g$	計算時間 (秒)
MPIM	40.907	0.35
SBMPIM	40.910±0.043	1.11
SMART	62.626	3.98
RELAXED-SMART	65.500	4.01
SBPI	43.217	1.92

3 品目在庫管理問題に対して MPIM と SBMPIM により計算を行った結果を表 5 に示す。ここで, 2 品目の場合と異なる SBMPIM のパラメータは  $x_0=(7,8,8), x^*=(5,6,6)$  である。また, 需要およびコストのパラメータは品目 1, 2 に関しては 2 品目の場合と同じであり, 品目 3 に関しては, 需要パラメータは品目 1 と同じ, コストパラメータは品目 2 と同じである。多品目での SBMPIM の優位性が示されている。

表 5 3 品目問題の  $g$  および計算時間

アルゴリズム	$g$	計算時間 (秒)
MPIM	54.890	45.37
SBMPIM	54.904±0.041	113.27

## 謝辞

本研究の一部は, 日本学術振興会平成 14-15 年度科学研究費補助金 (基盤研究 (B) (1), 課題番号 14380185) および (基盤研究 (C) (2), 課題番号 14580477) の助成を受けてなされたものである。

## 参考文献

- [1] R. A. Howard: *Dynamic Programming and Markov Processes*, MIT Press (1960) (関根, 羽島, 森共訳「ダイナミックプログラミングとマルコフ過程」, 培風館, 1971)
- [2] M. L. Puterman: *Markov Decision Processes*, John Wiley & Sons (1994)
- [3] K. Ohno, T. Ishigaki and T. Yoshii: "A New Algorithm for a Multi-item Periodic Review Inventory System", *ZOR-Math. Methods of Oper. Res.* Vol. 39, pp. 349-364, 1994.
- [4] S. Russell and P. Norvig, 古川康一 訳「エージェントアプローチ 人工知能」, 共立出版, 1997.
- [5] R. S. Sutton and A. G. Barto: *Reinforcement Learning*, MIT Press (1998) (三上, 皆川共訳「強化学習」, 森北出版, 2000)
- [6] D. P. Bertsekas and J. N. Tsitsiklis: *Neuro-Dynamic Programming*, Athena Scientific (1996)
- [7] R. V. Roy: "Neuro-dynamic programming: overview and recent trends," pp.431-459, in E. A. Feinberg and A. Schwartz ed. *Handbook of Markov Decision Processes*, Kluwer Academic Publishers (2002)
- [8] T. K. Das, A. Gosavi, S. Mahadevan and Nich. Marchallick: "Solving semi-Markov decision problem using average reward reinforcement learning", *Management Science*, Vol. 45, No.4, pp.560-574 (1999)
- [9] A. Gosavi: Doctor Thesis, <http://faculty.uscolo.edu/gosavi/thesis.html> (1999)
- [10] Y. He, M. C. Fu and S. I. Marcus: "A Simulation-based policy iteration algorithm for average cost unichain Markov decision processes", M. Laguna and J. L. G. Velarde ed, *Computing Tools for Modeling, Optimization and Simulation*, Kluwer Academic, pp.161-182 (2000)
- [11] A. Gosavi, N. Bandla and T. K. Das: "A reinforcement learning approach to a single leg airline revenue management problem with multiple fare classes and overbooking," *IIE Transactions*, Vol. 34, pp.729-742 (2002)
- [12] 大野勝久: "マルコフ決定過程", システムと制御, Vol. 29, No. 6, pp.333-341 (1985)
- [13] K. Ohno and K. Ichiki: "Computing optimal policies for controlled tandem queueing systems", *Operations Research*, Vol. 35, No. 1, pp.121-126 (1987)
- [14] K. Ohno: "Modified policy iteration algorithm with nonoptimality tests for undiscounted Markov decision process", Working Paper, Dept. of Information System and Management Science, Konan University, Japan (1985)
- [15] 大野勝久, 八嶋憲司, 伊藤崇博: "ニューロ・ダイナミックプログラミングによる生産ラインの最適制御に関する研究", 日本経営工学会論文誌, Vol. 54, No. 5, pp. 316-325 (2003)
- [16] E. L. Johnson: "Optimality and Computation of  $(\sigma, S)$  Policies in the Multi-item Infinite Horizon Inventory Problem", *Management Science*, Vol. 13, pp. 475-491 (1967)
- [17] 大野勝久, 田村隆善, 森健一, 中島健一: 「生産管理システム」, 朝倉書店 (2002)
- [18] 石塚陽, 山下英明: 「サンプルパス最適化の確率的離散事象システムへの適用」, オペレーションズ・リサーチ, Vol. 46, No. 4, pp.195-201 (2001)



表4 各アルゴリズムの政策の比較

状態		MPIM		SBMPIM			SBPI		
品目0	品目1	発注0	発注1	発注0	発注1	一致	発注0	発注1	一致
0	0	6	6	6	6	○	5	6	×
4	4	0	0	0	0	○	0	0	○
0	1	6	5	6	5	○	5	5	×
1	0	5	6	5	6	○	4	6	×
3	4	0	0	0	0	○	0	0	○
4	3	0	0	0	0	○	0	0	○
3	3	0	0	0	0	○	0	0	○
2	2	0	0	0	0	○	0	3	×
1	1	5	5	5	5	○	4	5	×
4	5	0	0	0	0	○	0	0	○
1	2	0	0	0	0	○	4	4	×
0	2	6	4	6	4	○	5	4	×
2	0	4	6	4	6	○	0	5	×
2	1	0	0	0	0	○	0	4	×
3	2	0	0	0	0	○	0	0	○
政策一致数 (300 状態中)				300			136		

状態		MPIM		SMART			RELAXED-SMART		
品目0	品目1	発注0	発注1	発注0	発注1	一致	発注0	発注1	一致
0	0	6	6	2	1	×	2	1	×
4	4	0	0	0	0	○	0	0	○
0	1	6	5	1	1	×	1	0	×
1	0	5	6	0	0	×	0	1	×
3	4	0	0	0	0	○	0	0	○
4	3	0	0	0	0	○	0	0	○
3	3	0	0	0	0	○	0	0	○
2	2	0	0	0	0	○	0	0	○
1	1	5	5	0	0	×	0	0	×
4	5	0	0	0	0	○	0	0	○
1	2	0	0	0	0	○	0	0	○
0	2	6	4	0	0	×	1	0	×
2	0	4	6	0	0	×	0	0	×
2	1	0	0	0	0	○	0	0	○
3	2	0	0	0	0	○	0	0	○
政策一致数 (300 状態中)				242			237		