

Maximizing Order Probabilities on Controlled Markov Chains

長崎県立大学経済学部 植野 貴之 (Takayuki Ueno)
Faculty of Economics, Nagasaki Prefectural University
九州大学大学院経済学研究院 岩本 誠一 (Seiichi Iwamoto)
Graduate School of Economics, Kyushu University

1 Introduction

It is natural to consider expectation criteria in stochastic decision problems. In particular, both discounted expected value of total reward and average value per stage are well studied criteria in Markov decision processes.

However, in some situation such as in economy of lower growth rate, growth of economy is in itself a preferable criterion to total amount of production such as the gross national production. How can we measure the growth of decision process? We take an order that sequence of earned rewards is nondecreasing in stage direction.

In this paper we consider a probability criterion that the reward is nondecreasing in time—an order probability—. We maximize the order probability on finite-horizon controlled Markov chains. We show that the policy class for maximization depends upon reward function's dependence on today's state, today's decision and tomorrow's state.

In Section 2, we formulate an optimization problem with order probability criterion on finite-stage controlled Markov chains. Section 3 derives a recursive equation for decision process where reward function is independent of today's decision. Section 4 considers process with reward function of today's state, today's decision and tomorrow's state. It is shown that a recursive relation is derived through imbedding method by expanding the original state space.

2 Decision Process with Order Probability

Throughout the paper, the following data is given :

$N \geq 2$ is an integer; the *total number of stages*

$X = \{s_1, s_2, \dots, s_l\}$ is a *finite state space*

$U = \{a_1, a_2, \dots, a_k\}$ is a *finite action space*

$r_n : X \times U \rightarrow R^1$ is an *n-th reward function* ($0 \leq n \leq N-1$)

$k : X \rightarrow R^1$ is a *terminal function*

$p = \{p(\cdot|\cdot, \cdot)\}$ is a *Markov transition law*

$$: p(y|x, u) \geq 0 \quad \forall (x, u, y) \in X \times U \times X, \quad \sum_{y \in X} p(y|x, u) = 1 \quad \forall (x, u) \in X \times U$$

$y \sim p(\cdot|x, u)$ denotes that next state y conditioned on state x and action u appears with probability $p(y|x, u)$.

Let an N -stage controlled Markov chain $\{(X_n, U_n)\}$ on finite state space X and finite decision space U be under a Markov transition law p . We maximize the order probability that the reward (random variable) will appear in ascent order

$$r_0(X_0, U_0) \leq r_1(X_1, U_1) \leq \cdots \leq r_{N-1}(X_{N-1}, U_{N-1}) \leq k(X_N).$$

The problem is to how to find an optimal policy which maximizes the order probability $P(r_0 \leq r_1 \leq \cdots \leq r_{N-1} \leq k)$. We focus our attention on policy class where the optimization should be taken.

The *order probability*

$$P(r_0(X_0, U_0) \leq r_1(X_1, U_1) \leq \cdots \leq r_{N-1}(X_{N-1}, U_{N-1}) \leq k(X_N)) \quad (1)$$

depends not only upon initial state but also upon when, where and what the decision maker will choose. We maximize the order probability in Markov class :

$$\begin{aligned} & \text{Maximize } P_{x_0}^\pi(r_0 \leq r_1 \leq \cdots \leq r_{N-1} \leq k_N) \\ P_0(x_0) \quad & \text{subject to (i)}_n \quad X_{n+1} \sim p(\cdot | x_n, u_n) \quad 1 \leq n \leq N \\ & \text{(ii)}_n \quad u_n \in U \end{aligned}$$

Here $P_{x_0}^\pi$ is the (discrete) probability measure on the product space X^N induced from the transition law p , a Markov policy $\pi = \{\pi_0, \pi_1, \dots, \pi_{N-1}\} \in \Pi$, and an initial state $x_0 \in X$. Thus the probability is expressed by the *partial* multiple summation :

$$P_{x_0}^\pi(r_0 \leq \cdots \leq k) = \sum_{(x_1, x_2, \dots, x_N) \in (*)} \sum \cdots \sum p(x_1 | x_0, u_0) p(x_2 | x_1, u_1) \cdots p(x_N | x_{N-1}, u_{N-1})$$

where the domain $(*)$ in which the partial multiple summation is taken denotes the set of all $(x_1, x_2, \dots, x_N) \in X \times X \times \cdots \times X$ satisfying

$$r_0(x_0, u_0) \leq r_1(x_1, u_1) \leq \cdots \leq r_{N-1}(x_{N-1}, u_{N-1}) \leq r(x_N).$$

Here we note that the sequence of intermediate decisions $\{u_0, u_1, \dots, u_{N-1}\}$ is determined through the Markov policy $\pi = \{\pi_0, \dots, \pi_{N-1}\}$ as follows :

$$u_0 = \pi_0(x_0), \quad u_1 = \pi_1(x_1), \quad \dots, \quad u_{N-1} = \pi_{N-1}(x_{N-1}).$$

Thus the order probability control problem is written as follows :

$$\begin{aligned} & \text{Maximize } P_{x_0}^\pi(r_0 \leq \cdots \leq k) \\ P_0(x_0) \quad & \text{subject to (i)}_n, \text{ (ii)}_n \quad 0 \leq n \leq N-1. \end{aligned}$$

We call this problem the *order problem*. Let $v_0(x_0)$ denote the *maximum value* of problem $P_0(x_0)$:

$$v_0(x_0) := \text{Max}_{\pi \in \Pi} P_{x_0}^\pi(r_0 \leq \cdots \leq k) \quad x_0 \in X. \quad (2)$$

Then our problem is to find an *optimal* policy π^* in Markov class Π :

$$v_0(x_0) = P_{x_0}^{\pi^*}(r_0 \leq \cdots \leq k) \quad \forall x_0 \in X. \quad (3)$$

3 Subproblems

In this section we derive a recursive equation for process $P_0(x_0)$. Let $\pi = \{\pi_n, \pi_{n+1}, \dots, \pi_{N-1}\}$ be a Markov policy for subprocess from n -th stage on. We denote by $\Pi(n)$ the set of all such Markov policies.

Now let us take n ($0 \leq n \leq N-1$), $x_n \in X$ and a policy $\pi \in \Pi(n)$. We consider the order probability

$$\begin{aligned} & P_{x_n}^\pi(r_n \leq r_{n+1} \leq \dots \leq k) \\ := & P_{x_n}^\pi(r_n(X_n, U_n) \leq r_{n+1}(X_{n+1}, U_{n+1}) \leq \dots \leq r_{N-1}(X_{N-1}, U_{N-1}) \leq k(X_N)). \end{aligned}$$

for the process which starts at state $x_n \in X$ and is governed by $\pi \in \Pi(n)$. Formally we set

$$P_{x_N}(k(X_N)) := 1.$$

Then we have the recursive relation:

Lemma 3.1 For $n \leq N-2$ we have

$$\begin{aligned} P_{x_n}^\pi(r_n \leq r_{n+1} \leq \dots \leq k) &= \sum_{*} P_{x_{n+1}}^{\pi'}(r_{n+1} \leq \dots \leq k) p(x_{n+1}|x_n, u_n) \\ &\hookrightarrow * : x_{n+1} \in X; r_n(x_n, u_n) \leq r_{n+1}(x_{n+1}, u_{n+1}) \end{aligned} \quad (4)$$

where $\pi' = \{\pi_{n+1}, \dots, \pi_{N-1}\}$, $u_m = \pi_m(x_m)$ $m = n, n+1$.

Further for $\pi = \{\pi_{N-1}\}$ we have

$$\begin{aligned} P_{x_{N-1}}^\pi(r_{N-1} \leq k) &= \sum_{*} P_{x_N}(k) p(x_N|x_{N-1}, u_{N-1}) \\ &\hookrightarrow * : x_N \in X; r_{N-1}(x_{N-1}, u_{N-1}) \leq k(x_N) \end{aligned} \quad (5)$$

where $u_{N-1} = \pi_{N-1}(x_{N-1})$.

We see that Eq. (5) states

$$P_{x_{N-1}}^\pi(r_{N-1} \leq k) = \sum_{*} p(x_N|x_{N-1}, u_{N-1})$$

Now we consider the family of subproblems:

$$\begin{aligned} & \text{maximize } P_{x_n}^\pi(r_n \leq r_{n+1} \leq \dots \leq r_{N-1} \leq k) \\ P_n(x_n) & \text{ subject to (i)}_m X_{m+1} \sim p(\cdot|x_m, u_m) \quad n \leq m \leq N \\ & \text{(ii)}_m u_m \in U \end{aligned}$$

where $0 \leq n \leq N-1$, $x_n \in X$. Let $v_n(x_n)$ denote the maximum value of problem $P_n(x_n)$:

$$v_n(x_n) := \text{Max}_{\pi \in \Pi(n)} P_{x_n}^\pi(r_n \leq \dots \leq k) \quad x_n \in X. \quad (6)$$

where we set

$$v_N(x_N) := 1 \quad x_N \in X. \quad (7)$$

Then we have the recursive equation:

Theorem 3.1

$$\begin{aligned}
v_N(x) &= 1 \quad x \in X \\
v_{N-1}(x) &= \text{Max}_{u \in U} \sum_{*} v_N(y) p(y|x, u) \quad x \in X \\
&\hookrightarrow * : y \in X; r_{N-1}(x, u) \leq k(y)
\end{aligned} \tag{8}$$

$$\begin{aligned}
v_n(x) &= \text{Max}_{u \in U} \sum_{*} v_{n+1}(y) p(y|x, u) \quad x \in X, \quad n = N-2, \dots, 1, 0. \\
&\hookrightarrow * : y \in X; r_n(x, u) \leq r_{n+1}^*(y)
\end{aligned} \tag{9}$$

where $\pi_{N-1}^*, \pi_{N-2}^*, \dots, \pi_1^*, \pi_0^*$ are calculated backward; the first $\pi_{N-1}^*(x)$ is a maximizer for (8) and the subsequent $\pi_n^*(x)$ is a maximizer for (9). Further, $r_{N-1}^*, r_{N-2}^*, \dots, r_1^*$ are successively defined through $\pi_{N-1}^*, \pi_{N-2}^*, \dots, \pi_1^*$:

$$r_n^*(x) := r_n(x, \pi_n^*(x)) \quad n = N-1, N-2, \dots, 1. \tag{10}$$

3.1 Decision-free Reward System

Now we consider the special case where the reward function $r : X \times U \rightarrow R^1$ is independent of decision variable n, u :

$$r(x, u) = r(x).$$

Thus the order probability is

$$P_{x_0}^\pi(r_0 \leq r_1 \leq \dots \leq k) := P_{x_0}^\pi(r(X_0) \leq r(X_1) \leq \dots \leq r(X_{N-1}) \leq k(X_N)). \tag{11}$$

Then we have simplified results.

Corollary 3.1 For $n \leq N-2$ we have

$$\begin{aligned}
P_{x_n}^\pi(r_n \leq r_{n+1} \leq \dots \leq k) &= \sum_{*} P_{x_{n+1}}^{\pi'}(r_{n+1} \leq \dots \leq k) p(x_{n+1}|x_n, u_n) \\
&\hookrightarrow * : x_{n+1} \in X; r(x_n) \leq r(x_{n+1})
\end{aligned} \tag{12}$$

where $u_n = \pi_n(x_n)$, $\pi' = \{\pi_{n+1}, \dots, \pi_{N-1}\}$.

Further for $\pi = \{\pi_{N-1}\}$ we have

$$\begin{aligned}
P_{x_{N-1}}^\pi(r_{N-1} \leq k) &= \sum_{*} P_{x_N}(k) p(x_N|x_{N-1}, u_{N-1}) \\
&\hookrightarrow * : x_N \in X; r(x_{N-1}) \leq k(x_N)
\end{aligned} \tag{13}$$

where $u_{N-1} = \pi_{N-1}(x_{N-1})$.

Corollary 3.2

$$\begin{aligned}
v_N(x) &= 1 \quad x \in X \\
v_{N-1}(x) &= \text{Max}_{u \in U} \sum_{\substack{y \in X \\ r(x) \leq k(y)}} v_N(y) p(y|x, u) \quad x \in X \\
v_n(x) &= \text{Max}_{u \in U} \sum_{\substack{y \in X \\ r(x) \leq r(y)}} v_{n+1}(y) p(y|x, u) \quad x \in X, \quad n = N-2, \dots, 0.
\end{aligned} \tag{14}$$

4 Reward Functions Depend on Tomorrow

In this section we treat the reward function which depends not only today but also on tomorrow. We consider both evaluation problem of order probability and optimization problem.

4.1 Evaluation

First we consider a recursive evaluation of order probability. Now let a sequence of reward functions

$$r_n : X \times X \rightarrow R^1 \quad (0 \leq n \leq N-1), \quad k : X \rightarrow R^1$$

be given. We note that the reward functions depend on next state.

$$P_0(x_0) \quad \begin{array}{l} \text{Evaluate } P_{x_0}(r_0 \leq r_1 \leq \cdots \leq r_{N-1} \leq k) \\ \text{under (i)}_n \quad X_{n+1} \sim p(\cdot | x_n) \quad n = 0, \dots, N-1 \end{array}$$

Thus we evaluate the order probability

$$v_0(x_0) := P_{x_0}(r_0 \leq r_1 \leq \cdots \leq r_{N-1} \leq k) \quad (15)$$

under the Markov chain $\{X_n\}_0^N$ with transition probability law $p = \{p(\cdot | \cdot)\}$.

Let us consider the family of subproblems $\{P_n(x_n, x_{n+1})\}$:

$$P_n(x_n, x_{n+1}) \quad \begin{array}{l} \text{Evaluate } P_{x_n, x_{n+1}}(r_n \leq r_{n+1} \leq \cdots \leq r_{N-1} \leq k) \\ \text{under (i)}_m \quad X_{m+1} \sim p(\cdot | x_m) \quad m = n+1, \dots, N-1 \\ (x_n, x_{n+1}) \in X \times X, \quad n = 0, \dots, N-1. \end{array}$$

Here we note that the evaluated order probability is the conditional probability :

$$\begin{aligned} & P_{x_n, x_{n+1}}(r_n \leq r_{n+1} \leq \cdots \leq r_{N-1} \leq k) \\ = & P(r_n(X_n, X_{n+1}) \leq r_{n+1}(X_{n+1}, X_{n+2}) \leq \cdots \leq r_{N-1}(X_{N-1}, X_N) \leq k(X_N) \\ & \quad | (X_n, X_{n+1}) = (x_n, x_{n+1})) \\ = & \sum_{* ; x_{n+2}, \dots, x_N} \cdots \sum p(x_{n+2} | x_{n+1}) p(x_{n+3} | x_{n+2}) \cdots p(x_N | x_{N-1}, x_N) \\ \hookrightarrow & * : r_n(x_n, x_{n+1}) \leq r_{n+1}(x_{n+1}, x_{n+2}) \leq \cdots \leq r_{N-1}(x_{N-1}, x_N) \leq k(x_N). \end{aligned}$$

Let $w_n(x_n, x_{n+1})$ denote the probability of $P_n(x_n, x_{n+1})$, where

$$w_N(x_N) := 1.$$

Then we have the recursive equation

Lemma 4.1

$$\begin{aligned} w_N(x) &= 1 & x \in X \\ w_{N-1}(x, y) &= \begin{cases} 1 & \text{if } r_{N-1}(x, y) \leq k(y) \\ 0 & \text{otherwise} \end{cases} & (x, y) \in X \times X \\ w_n(x, y) &= \sum_{*} w_{n+1}(y, z) p(z | y) & (x, y) \in X \times X, \quad 0 \leq n \leq N-1 \\ \hookrightarrow & * : z \in X ; r_n(x, y) \leq r_{n+1}(y, z) \end{aligned} \quad (16)$$

The desired probability (15) is given by

$$v_0(x_0) = \sum_{x_1 \in X} w_0(x_0, x_1) p(x_1 | x_0)$$

4.2 Optimization

Second we consider a recursive optimization of order probability. Let reward functions

$$r_n : X \times U \times X \rightarrow R^1 \quad (0 \leq n \leq N-1), \quad k : X \rightarrow R^1$$

be dependent on well next state as current decision. We optimize the order probability

$$P_{x_0}(r_0 \leq r_1 \leq \dots \leq r_{N-1} \leq k) \quad (17)$$

under the controlled Markov chain $\{(X_n, U_n)\}$ with transition probability law $p = \{p(\cdot | \cdot, \cdot)\}$. The problem is which class we optimize in and how we can find an optimal policy. The preceding discussion on evaluation enables us to choose an policy class where any policy is a sequence of decision functions :

$$\gamma_n : X \times \Lambda_n \rightarrow U \quad 2 \leq n \leq N-1.$$

Let us introduce the sequence of *yesterday (last) reward sets* $\{\Lambda_n(x_n)\}$ to current state x_n :

$$\begin{aligned} \Lambda_0(x_0) &\triangleq \{\lambda_0 \mid \lambda_0 = -\infty\} \\ \Lambda_n(x_n) &\triangleq \left\{ \lambda_n \mid \begin{array}{l} \lambda_n = r_{n-1}(x_{n-1}, u_{n-1}, x_n) \\ (x_{n-1}, u_{n-1}) \in X \times U \end{array} \right\} \\ &\quad x_n \in X, \quad n = 1, \dots, N. \end{aligned} \quad (18)$$

Further we define *yesterday (last) reward set*

$$\Lambda_n := \bigcup_{x \in X} \Lambda_n(x)$$

and the sequence of *expanded state spaces* $\{Y_n\}$:

$$Y_n := \{y_n = (x_n, \lambda_n) \mid x_n \in X, \lambda_n \in \Lambda_n(x_n)\} \quad n = 0, \dots, N.$$

Let us define the corresponding *random variable* $\tilde{\Lambda}_n$ by

$$\begin{aligned} \tilde{\Lambda}_0 &\triangleq \lambda_0 \\ \tilde{\Lambda}_n &\triangleq r_{n-1}(X_{n-1}, U_{n-1}, X_n) \quad n = 1, \dots, N. \end{aligned}$$

which takes values in Λ_n .

Now we introduce a new controlled Markov chain on the expanded state spaces $\{Y_n\}_0^N$. Here the state variables $\{(X_n; \tilde{\Lambda}_n)\}$ behave such that the first component $\{X_n\}$ obeys the original Markov transition law p and the second $\{\tilde{\Lambda}_n\}$ follows the stochastic dynamics $\tilde{\Lambda}_{n+1} := r_n(x_n, u_n, X_{n+1})$. When the decision-maker chooses a decision $u_n (\in U)$ on $(x_n; \lambda_n) (\in Y_n)$ at n -th stage, the next state random variable $(X_{n+1}; \tilde{\Lambda}_{n+1})$ will take $(x_{n+1}; \lambda_{n+1})$ with probability

$p(x_{n+1}|x_n, u_n)$ at $(n+1)$ -st stage, where $\lambda_{n+1} = r_n(x_n, u_n, x_{n+1})$. Thus this is expressed by a coupled dynamics (i)_n, (i')_n $0 \leq n \leq N-1$.

Thus we maximize the order probability over the new controlled chain on expanded state spaces as follows :

$$\begin{aligned} & \text{maximize } P_{x_0}^\gamma(\lambda_0 \leq r_0 \leq r_1 \leq \dots \leq r_{N-1} \leq k) \\ P_0(x_0; \lambda_0) & \text{ subject to (i)}_n \quad X_{n+1} \sim p(\cdot | x_n, u_n) \\ & \text{(i')}'_n \quad \tilde{\Lambda}_{n+1} = r_n(x_n, u_n, X_{n+1}) \quad 0 \leq n \leq N-1 \\ & \text{(ii)}_n \quad u_n \in U \end{aligned}$$

where $\lambda_0 = -\infty$.

We imbed $P_0(x_0; \lambda_0)$ into the family of subproblems $\{P_n(x_n; \lambda_n)\}$:

$$\begin{aligned} & \text{maximize } P_{x_n}^\gamma(\lambda_n \leq r_n \leq r_{n+1} \leq \dots \leq r_{N-1} \leq k) \\ P_n(x_n; \lambda_n) & \text{ subject to (i)}_m \quad X_{m+1} \sim p(\cdot | x_m, u_m) \\ & \text{(i')}'_m \quad \tilde{\Lambda}_{m+1} = r_m(x_m, u_m, X_{m+1}) \quad n \leq m \leq N-1 \\ & \text{(ii)}_m \quad u_m \in U \end{aligned}$$

where $x_n \in X, \lambda_n \in \Lambda_n(x_n)$ and $0 \leq n \leq N-1$.

Let $v_n(x_n, \lambda_n)$ denote the maximum value of $P_n(x_n; \lambda_n)$, where we set

$$v_N(x_N, \lambda_N) := P_{x_N}(\lambda_N \leq k(x_N)).$$

Then we have the recursive equation:

Theorem 4.1

$$v_N(x; \lambda) = \begin{cases} 1 & \text{if } \lambda \leq k(x) \\ 0 & \text{otherwise} \end{cases} \quad x \in X, \lambda \in \Lambda_N(x)$$

$$v_{N-1}(x; \lambda) = \text{Max}_{u \in U} \sum_{\substack{* \\ \hookrightarrow *}} v_N(y; r_{N-1}(x, u, y)) p(y|x, u) \quad x \in X, \lambda \in \Lambda_{N-1}(x) \quad (19)$$

$$v_n(x; \lambda) = \text{Max}_{u \in U} \sum_{\substack{* \\ \hookrightarrow *}} v_{n+1}(y; r_n(x, u, y)) p(y|x, u) \quad x \in X, \lambda \in \Lambda_n(x) \quad n = N-2, \dots, 0 \quad (20)$$

where $\pi_{N-1}^*, \pi_{N-2}^*, \dots, \pi_1^*, \pi_0^*$ are calculated backward ; the first $\pi_{N-1}^*(x; \lambda)$ is a maximizer for (19) and the subsequent $\pi_n^*(x; \lambda)$ is a maximizer for (20).

References

- [1] Altman, E.: Constrained Markov Decision Processes. Chapman & Hall, New York, 1999
- [2] Bellman, R.E.: Dynamic Programming. Princeton University Press, Princeton, NJ, 1957

- [3] Bellman, R.E.: Some Vistas of Modern Mathematics. University of Kentucky Press, Lexington, KY, 1968
- [4] Blackwell, D.: Discounted dynamic programming. *Ann. Math. Stat.* **36**, 226-235 (1965)
- [5] Denardo, E.V.: Contraction mappings in the theory underlying dynamic programming. *SIAM Review* **9**, 165-177 (1968)
- [6] Denardo, E.V.: *Dynamic Programming: Models and Applications*. Prentice-Hall, NJ, 1982
- [7] Dynkin, E.B., Yushkevich, A.A.: *Controlled Markov Processes*. Springer, New York, 1979
- [8] Hinderer, K.: *Foundations of Non-Stationary Dynamic Programming with Discrete Time Parameter*. *Lectures Notes in Operation Research and Mathematical Systems* **33**, Springer, Berlin, 1970
- [9] Howard, R.A.: *Dynamic Programming and Markov Processes*. MIT Press, Cambridge, Mass., 1960
- [10] Iwamoto, S.: *Theory of Dynamic Program (in Japanese)*. Kyushu Univ. Press, Fukuoka, 1987
- [11] Iwamoto, S.: Associative dynamic programs. *J. Math. Anal. Appl.* **201**, 195-211 (1996)
- [12] Iwamoto, S.: On expected values of Markov statistics. *Bull. Informatics and Cybernetics* **30**, 1-24 (1998)
- [13] Iwamoto, S.: Conditional decision processes with recursive reward function. *J. Math. Anal. Appl.* **230**, 193-210 (1999)
- [14] Iwamoto, S.: "Dynamic Programming", "Principle of Invariant Imbedding" (Japanese) In: *Operations Res. Soc. Japan (ed.): Operations Research Dictionary 2000: Basic Ver.*, pp.229-245, & *Terminology Ver.* JUSE, Tokyo, 2000
- [15] Iwamoto, S., Fujita, T.: Stochastic decision-making in a fuzzy environment. *J. Operations Res. Soc. Japan* **38**, 467-482 (1995)
- [16] Iwamoto, S., Sniedovich, M.: Sequential decision making in fuzzy environment. *J. Math. Anal. Appl.* **222**, 208-224 (1998)
- [17] Iwamoto, S., Tsurusaki, K., Fujita, T.: Conditional decision-making in a fuzzy environment. *J. Operations Res. Soc. Japan* **42**, 198-218 (1999)
- [18] Iwamoto, S., Ueno, T., Fujita, T.: Controlled Markov chains with utility functions. *International Workshop on Markov Processes and Controlled Markov Chains, Changsha, Hunan, China, August 22-28, 1999*
- [19] Karatzas, I., Shreve, S.E.: *Methods of Mathematical Finance*. Springer, New York, 1998.
- [20] Kreps, D.M.: Decision problems with expected utility criteria I. *Math. Oper. Res.* **2**, 45-53 (1977)

- [21] Kreps, D.M.: Decision problems with expected utility criteria II; stationarity. *Math. Oper. Res.* **2**, 266-274 (1977)
- [22] Markovitz, H.: Portfolio selection. *J. Finance* **8**, 77-91 (1952)
- [23] Ozaki, H., Streufert, P.A.: Dynamic programming for non-additive stochastic objects. *J. Math. Eco.* **25**, 391-442 (1996)
- [24] Porteus, E.: An informal look at the principle of optimality. *Management Sci.* **21**, 1346-1348 (1975)
- [25] Porteus, E.: Conditions for characterizing the structure of optimal strategies in infinite-horizon dynamic programs. *J. Opt. Theo. Anal.* **36**, 419-432 (1982)
- [26] Puterman, M.L.: Markov Decision Processes: Stochastic Models. In: Heyman, D.P., Sobel, M.J. (eds.): *Handbooks in Operations Research and Management Science Vol. 2*, Elsevier, Amsterdam, 1990, Chap. VIII
- [27] Puterman, M.L.: *Markov Decision Processes: discrete stochastic dynamic programming*. Wiley & Sons, New York, 1994
- [28] Sniedovich, M.: *Dynamic Programming*. Marcel Dekker, Inc. NY, 1992
- [29] Stokey, N.L., Lucas, R.E.: *Recursive Methods in Economic Dynamics*. Harvard University Press, Cambridge, Mass., 1989
- [30] Strauch, R.: Negative dynamic programming. *Ann. Math. Stat.* **37**, 871-890 (1966)
- [31] Streufert, P.A.: Recursive Utility and Dynamic Programming. In: Barberà, S. et al. (eds.): *Handbook of Utility Theory Vol. 1*, Kluwer, Boston, 1998, Chap. III