

# Improved successive approximation methods for discounted Markov decision processes

Citation for published version (APA): van Nunen, J. A. E. E. (1974). Improved successive approximation methods for discounted Markov decision processes. (Memorandum COSOR; Vol. 7406). Technische Hogeschool Eindhoven.

## Document status and date:

Published: 01/01/1974

#### Document Version:

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

#### Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

Link to publication

#### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- · Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
  You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

#### Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

Download date: 22. Aug. 2022

# EINDHOVEN UNIVERSITY OF TECHNOLOGY Department of Mathematics

STATISTICS AND OPERATION RESEARCH GROUP

Memorandum COSOR 74-06

Improved successive approximation methods for discounted Markov decision processes

bу

J.A.E.E. van Nunen

#### Abstract

Successive Approximation (S.A.) methods, for solving discounted Markov decision problems, have been developed to avoid the extensive computations that are connected with linear programming and policy iteration techniques for solving large scaled problems. Several authors give such an S.A. algorithm.

In this paper we introduce some new algorithms while furthermore it will be shown how the several S.A. algorithms may be combined. For each algorithm converging sequences of upper and lower bounds for the optimal value will be given.

### § 1. Introduction.

We consider a finite state, discrete time Markov system that is controlled by a decisionmaker (see for example [4]). After each transition n = 0,1,2,... the system may be identified as being in one of N possible states. Let  $S := \{1,2,\ldots,N\}$  represent the set of states. After observing state  $i \in S$  the decisionmaker selects an action k from a nonempty finite set K(i). Now  $p_{ij}^k$  is the probability of a transition to state  $j \in S$ , if the system is actually in state  $i \in S$ , and action  $k \in K(i)$  has been selected. An (expected) reward  $q^k(i)$  is then earned immediately, while future income is discounted by a constant factor  $0 \le \alpha < 1$ .

We suppose, which is permitted without loss of generality, that  $q^k(i) \ge 0$  for all  $i \in S$  and  $k \in K(i)$ .

The problem is to choose a policy which maximizes the total expected discounted return.

As known (e.g. [2], [10]), it is permitted to restrict the considerations to nonrandomized stationary strategies. A nonrandomized stationary strategy will be denoted by  $f \in K := K(1) \times K(2) \times \ldots \times K(N)$ . The coordinates  $u_f(i)$  of the N  $\times$  1 vector  $u_f$  give the total expected discounted return if the system's initial state is i and the stationary strategy  $f \in K$  is used. The (stationary) strategy  $f \in K$  is called optimal if  $u_f(i) \ge u_f(i)$  for all  $f \in K$  and for all  $i \in S$ .

Because S.A. algorithms are in some sense modifications of the standard dynamic programming method this method will be discussed first.

As in Blackwell [1] we define for each  $f \in K$  the mapping  $L_b(f) (\mathbb{R}^N \to \mathbb{R}^N)$  which maps an N × 1 column vector x into

$$L_h(f)x := q_f + \alpha P_f x$$
,

where  $q_f$  is the N × 1 column vector having as its i-th component  $q^{f(i)}(i)$ , and  $P_f$  is the N × N Markov matrix with (i,j) element  $p_{ij}^{f(i)}$ .  $L_{b}^{(f)}(f)$  is monotone, i.e. if every coordinate of the N × 1 vector x is at least as large as every coordinate of  $y \in \mathbb{R}^N$   $(x \ge y)$ , then:

$$L_b(f)x \ge L_b(f)y$$
.

Furthermore, we define for some map  $L_{\beta}(f)$ :

$$L_{\beta}^{0}(f)x := x;$$

$$L_{\beta}^{n}(f)x = L_{\beta}(f)(L_{\beta}^{n-1}(f)x) .$$

We define the mapping  $U_b: \mathbb{R}^N \to \mathbb{R}^N$  by:

$$U_b x := \max_{f \in K} L_b(f) x$$
.

It is easily seen that for every  $x \in \mathbb{R}^N$  an  $f \in K$  exists such that  $L_b(f)x$  is maximal for each coordinate.

It may be proved that  $U_b$  is a monotone  $\alpha$ -contraction mapping with fixed point  $u_{\star}$ . For an optimal strategy  $f_{\star}$  we have  $u_{\star} = u_{f_{\star}}$ , and a stationary strategy f which is optimal follows from  $L_b(f)u_{\star} := U_bu_{\star}$  (see for instance [9]).

This property legitimates the standard dynamic programming algorithm that can be based on:

$$I \begin{cases} x_0^b := 0 \\ x_n^b := U_b x_{n-1}^b =: L_b(f_n^b) x_{n-1}^b \end{cases}$$

It is possible to take  $x_n^b$  and  $f_n^b$  as estimates for  $u_{f_{\star}}$  and  $f_{\star}$  as follows from:

(1) 
$$x_{n-1}^{b} \leq x_{n}^{b} \leq u_{f_{n}}^{b} \leq u_{f_{\star}}^{\star}$$

(2) 
$$\lim_{n\to\infty} x_n^b = \lim_{n\to\infty} U_b^n x_0^b = u_{f_*},$$

see [5], [7], [9].

As starting vector we choose  $x_0^b = 0$ . As appears from (1) and (2), this choice guarantees monotone convergence of  $x_n^b$  to  $u_f$ .

As known, the convergence, depending on  $\alpha$ , may be relatively slow. MacQueen [5] constructed upper and more sophisticated lower bounds for u and uf.

The S.A. methods discussed in the following sections are based on contraction mappings (see Denardo [2]).

In section 2, S.A. methods based on mappings  $U_0$ ,  $U_a$ ,  $U_s$ ) of the same type as  $U_b$  will be given; i.e. also these mappings are monotone contraction mappings with fixed point  $u_{f\star}$ . Furthermore, combinations of these mappings which lead to mappings  $(U_{hs}, U_{h0})$  with the same property will be discussed. In section 3 extension of the above algorithms will be given, while in section 4 upper and lower bounds for several methods are discussed. This enables us to incorporate a test for the suboptimality of actions (see also [6]). Finally (section 5) some examples are given to illustrate several methods.

# § 2. "Improved" successive approximation methods.

2.1. Hastings [3] introduced the following (Gauss-Seidel) idea to modify the policy improvement procedure in Howard's policy iteration algorithm. Let  $u_f$  for a given strategy  $f \in K$  be computed. Determine a better strategy  $g \in K$  with components g(i) as follows:

g(1) follows from 
$$v_g(1) := \max_{k \in K(1)} \{q^k(1) + \alpha \sum_{j=1}^{N} p_{ij}^k u_f(j)\}$$
  

$$=: \{q^{g(1)}(1) + \alpha \sum_{j=1}^{N} p_{ij}^{g(1)} u_f(j)\},$$

This idea can also be used in an S.A. algorithm.

Let  $x \in \mathbb{R}^N$ . Define  $L_h(f)$  by:

$$L_h(f)x(i) := q^{f(i)}(i) + \alpha \sum_{j < i} p_{ij}^{f(i)}(L_h(f)x(j)) + \alpha \sum_{j \ge i} p_{ij}^{f(i)}x(j) , \quad i \in S.$$

Define the mapping U<sub>h</sub> by:

$$U_h x := \max_{f} L_h(f)x$$
.

It is easily verified that  $L_h(f)$  and  $U_h$  are monotone  $\alpha$ -contractions with fixed point  $u_f$  and  $u_{f_+}$ , respectively, so an S.A. algorithm might be based on

II 
$$\begin{cases} x_0^h := 0 \\ x_n^h := U_h x_{n-1}^h =: L_h(f_n^h) x_{n-1}^h \end{cases}$$

As in standard dynamic programming, the sequence  $\{x_n^h\}$  will have the following properties:

(3) 
$$x_{n-1}^{h} \le x_{n}^{h} \le u_{f_{n}}^{h} \le u_{f_{*}}^{h}$$
,

(4) 
$$\lim_{n\to\infty} x_n^h := u_{f_*}.$$

Furthermore, a comparison with the  $x_n^b$  of the dynamic programming algorithm yields inductively

$$(5) x_n^h \ge x_n^b.$$

2.2. Also "overrelaxation" (see [8], [9]) may be used in successive approximation algorithms. Where the overrelaxation factor appears, for instance, if we try to find better estimates for uf by computing for certain paths the exact contribution to the total expected discounted reward.

Let  $f \in K$  be given, then

$$u_f = L_{f}u_f = q_f + \alpha P_{f}u_f$$
.

Another expression for  $u_f(i)$  may be found by computing the contribution to the expected reward until the time the system leaves i explicitly (f(i) =: k):

(6) 
$$u_{f}(i) = q^{k}(i) + \alpha p_{ii}^{k} q^{k}(i) + (\alpha p_{ii}^{k})^{2} q^{k}(i) + \dots$$

$$+ \alpha \sum_{j \neq i} p_{ij}^{k} u_{f}(j)$$

$$+ \alpha^{2} p_{ii}^{k} \sum_{j \neq i} p_{ij}^{k} u_{f}(j)$$

$$\vdots$$

$$= \frac{1}{1 - \alpha p_{ii}^{k}} q^{k}(i) + \frac{\alpha}{1 - \alpha p_{ii}^{k}} \sum_{j \neq i} p_{ij}^{k} u_{f}(j) .$$

Let  $\omega_i^k := \frac{1}{1 - \alpha p_{ij}^k}$ , then with k = f(i), (6) can be given as

(7) 
$$u_{f}(i) = \omega_{i}^{k} q^{k}(i) + \alpha \omega_{i}^{k} \sum_{j \neq i} p_{ij}^{k} u_{f}(j)$$
.

(7) can also be deduced from:

$$u_f(i) = q^k(i) + \alpha \sum_{j} p_{ij}^k u_f(j)$$
,

which yields:

$$(1 - \alpha p_{ii}^k) u_f(i) = q^k(i) + \alpha \sum_{i \neq i} p_{ij}^k u_f(j)$$
,

where (7) follows by dividing by  $(1 - \alpha p_{ii}^k)$ .

On the idea used in (7) we base for any  $f \in K$  the mapping  $L_0(f)$  defined by:

$$L_0(f)x(i) := \omega_i^k q^k(i) + \alpha \omega_i^k \sum_{j \neq i} p_{ij}^k x(j) \quad \text{with } k = f(i) .$$

Furthermore we define the mapping U<sub>0</sub> by:

$$U_0^x = \max_{f \in K} \{L_0(f)_x\}$$
.

Let

$$\omega^{-}(f) := \min_{i \in S} \{\omega_i^{f(i)}\}$$
,

$$\omega^+(f) := \max_{i \in S} \{\omega_i^{f(i)}\}$$
,

$$\gamma(\omega) := 1 - \omega(1 - \alpha) .$$

Then  $L_0(f)$  is a monotone  $\gamma(\omega^-(f))$ -contraction with fixed point  $u_f$ . It is easily verified that  $\gamma(\omega^-(f)) \leq \alpha$ .

Let  $\omega^*:=\min_{i,k}\{\omega_i^k\}$ , then  $U_0$  is a monotone  $\gamma(\omega^*)$ -contraction with fixed point  $u_{f_*}$  (see [8]).

We have the relation:

$$\gamma(\omega^{-}(f)) \leq \gamma(\omega^{*}) \leq \alpha$$
.

Hence a successive approximation method might be based on

III 
$$\begin{cases} x_0^0 := 0 \\ x_n^0 := U_0 x_{n-1}^0 =: L_0(f_n^0) x_{n-1}^0 \end{cases} ,$$

where the following inequalities are easily proved:

(8) 
$$x_{n-1}^{0} \le x_{n}^{0} \le u_{f_{n}}^{0} \le u_{f_{*}}^{0},$$

$$(9) x_n^b \le x_n^0.$$

2.4. It is also possible to simplify algorithm III by using the fixed overrelaxation factor  $\omega^*$ , which means that the contribution to the expected reward until the system leaves state i is only estimated.

Then we define  $L_s(f)$  by:

$$L_{s}(f)x(i) := \omega^{*} q^{k}(i) + \alpha\omega^{*} \sum_{j \in S} p_{ij}^{k} x(j) + (1 - \omega^{*})x(i) .$$

Ug is defined by:

$$U_{s}x := \max_{f \in K} L_{s}(f)x$$
.

 $L_s(f)$  and  $U_s$  are monotone  $\gamma(\omega^*)$ -contractions with fixed point  $u_f$  and  $u_{f_*}$ , respectively. So it is possible to construct an S.A. algorithm based on:

IV 
$$\begin{cases} x_0^s := 0 \\ x_n^s := U_s x_{n-1}^s =: L_s(f_n^s) x_{n-1}^s \end{cases}$$

Again we have:

(10) 
$$x_{n-1}^{s} \leq x_{n}^{s} \leq u_{f_{n}}^{s} \leq u_{f_{\star}}^{s}$$

$$(11) x_n^b \le x_n^s.$$

## § 3. Combinations of S.A. algorithms.

In this section it will be shown that combinations of the mappings  $\mathbf{U}_{h}$ ,  $\mathbf{U}_{0}$ ,  $\mathbf{U}_{s}$  lead to mappings  $\mathbf{U}_{h0}$ ,  $\mathbf{U}_{hs}$ , with the same properties as the original mappings; i.e.  $\mathbf{U}_{h0}$ ,  $\mathbf{U}_{hs}$  are monotone contractions with fixed point  $\mathbf{u}_{f}$ .

First we want to combine the transformations  $\mathbf{U}_0$  and  $\mathbf{U}_h$  as is done in a modified form by Reetz [8].

We define the transformation  $L_{h0}(f)$  inductively by

and U<sub>hO</sub> by:

$$U_{h0}x = \max_{f} L_{h0}(f)x$$
.

Then,  $L_{h0}(f)$  and  $U_{h0}$  are monotone and  $\gamma(\omega^*)$ -contractions with fixed point  $u_f$  and  $u_{f_*}$ , respectively.

So we have

$$V \begin{cases} x_0^{h0} = 0 \\ x_n^{h0} := U_{h0} x_{n-1}^{h0} =: L_{h0} (f_n^{h0}) x_{n-1}^{h0} \end{cases}$$

with

(12) 
$$x_{n-1}^{h0} \le x_n^{h0} \le u_{f_n}^{h0} \le u_{f_*},$$

$$\lim_{n \to \infty} x_n^{h0} = u_{f_*}.$$

Furthermore,

$$\mathbf{x}_{n}^{0} \leq \mathbf{x}_{n}^{h0} ,$$

$$(14) x_n^h \le x_n^{h0} .$$

The original Reetz [8] algorithm can be found as a combination of the transformations  $\mathbf{U_s}$  and  $\mathbf{U_h}$ .

Let  $L_{hs}(f)$  be given by

$$\begin{split} L_{hs}(f)x(i) &:= \omega^* q^k(i) + \alpha \omega^* \sum_{j < i} p_{ij}^k \ L_{hs}(f)x(j) \\ &+ \alpha \omega^* \sum_{j \ge i} p_{ij}^k \ x(j) + (1 - \omega^*)x(i) \end{split}$$

and

$$U_{hs}(f)x = \max_{f} L_{hs}(f)x$$
.

 $L_{\rm hs}(f)$  and  $U_{\rm hs}$  are monotone and  $\gamma(\omega^*)\text{-contractions}$  with fixed point  $u_f$  and  $u_f$  , respectively. We have:

$$VI \begin{cases} x_0^{hs} = 0 \\ x_n^{hs} = U_{hs}^{hs} x_{n-1}^{hs} =: L_{hs}(f_n^{hs}) x_{n-1}^{hs} \end{cases}$$

with

$$\begin{aligned} x_{n-1}^{hs} &\leq x_n^{hs} &\leq u_{f_n}^{hs} &\leq u_{f_n}^{}, \\ \lim_{n \to \infty} x_n^{hs} &= u_{f_n}^{}, \\ x_n^{s} &\leq x_n^{hs}, \\ x_n^{h} &\leq x_n^{hs}. \end{aligned}$$

# § 4. Extensions of S.A. algorithms.

A method to improve the estimations for  $u_{f_{\star}}$  can also be found by inserting a number of value determination iteration steps in the S.A. algorithm based on  $U_{\beta}$  where  $\beta \in T := \{b,h,0,s,hs,h0\}$ , see [7].

This idea can also be introduced as the skipping of a number of policy improvement iteration steps in the S.A. algorithms.

We define for each  $x\in {\rm I\!R}^N$  and for finite  $\lambda\in {\rm I\!R}$  and for  $\beta\in T$  the mapping  $U_\beta^{(\lambda)}$  by:

$$U_{\beta}^{(\lambda)} \mathbf{x} := L_{\beta}^{\lambda-1}(\mathbf{f}^{\beta\lambda}) U_{\beta} \mathbf{x}$$

where  $f^{\beta\lambda}$  indicates the strategy that is found by applying  $\textbf{U}_{\textrm{R}}$  on x.

For  $\lambda \in \mathbb{N}$ ,  $\lambda > 1$ ,  $U_{\beta}^{(\lambda)}$  is neither necessarily a contraction mapping nor a monotone mapping.

However, we may base an algorithm on such a mapping:

$$\begin{aligned} \text{VII-XII} \; & \left\{ \begin{array}{l} \mathbf{x}_0^{\beta\lambda} = 0 \;\; , & \beta \in T \\ \\ \mathbf{x}_n^{\beta\lambda} = \; \mathbf{U}_\beta^{(\lambda)} \mathbf{x}_{n-1}^{\beta\lambda} \; = : \; \mathbf{L}_\beta^\lambda(\mathbf{f}_n^{\beta\lambda}) \mathbf{x}_{n-1}^{\beta\lambda} \;\; , & \beta \in T \;\; . \end{array} \right. \\ \end{aligned}$$

The monotone convergence of  $x_n^{\beta\lambda}$  to  $u_{f_{\star}}$  is preserved (see [7]) as follows from the monotonicity of  $U_{g}$  and  $L_{g}(f^{\beta})$ , i.e.

$$\lim_{n\to\infty} x_n^{\beta\lambda} = u_{f_{\star}},$$

$$x_{n-1}^{\beta\lambda} \le x_n^{\beta\lambda} \le u_{f_n}^{\beta\lambda} \le u_{f_*}$$
.

A comparison of  $x_n^{\beta\lambda}$  with  $x_n^\beta$  yields:

$$\mathbf{x}_n^{\beta} \leq \mathbf{x}_n^{\beta \lambda}$$
 ,  $n \in \mathbb{N}, \ \beta \in \mathbb{T}$  .

# § 5. Upper and lower bounds for $u_{f_{\star}}$ .

Successive approximation algorithms based on the ideas of the previous sections will converge. However, it will be necessary to construct upper and lower bounds for the current and the optimal strategy. Upper and lower bounds enable us to qualify the estimates for  $u_{f_n}$ ,  $u_{f_n}$  and  $f_{\star}$ , respectively, see for instance MacQueen [5].

Also upper and lower bounds enable us to incorporate a test for the suboptimality of decisions in an algorithm (see MacQueen [6]).

Let the upper bound  $\bar{x}$  and the lower bound  $\underline{x}$  for  $u_{\star}$  be given, then we can state the following lemma:

Lemma 1. Strategy f is suboptimal if

$$L_{\beta}(f)\bar{x} < U_{\beta} \underline{x}$$
,  $\beta \in T$ .

Proof.

$$u_{\star} = U_{\beta}u_{\star} \ge U_{\beta}\underline{x} > L_{\beta}(f)\overline{x} \ge L_{\beta}(f)u_{\star}$$
,

where the monotonicity property of  $U_{\beta}$  and  $L_{\beta}(f)$  is used.

This lemma enables us to determine for each i  $\epsilon$  S decisions which are sub-optimal (see for instance [6]). \*

П

<sup>\*</sup> Note that in the algorithms where U  $_{\rm S}$  is used,  $\omega^*$  can be redefined if the decision that causes  $\omega^*$  is suboptimal.

If we want to compare two algorithms it will be necessary to compare the corresponding sequences of upper and lower bounds. However, where the estimates for  $\mathbf{u_f}$  found in the n-th iteration step of a specific algorithm may be better than those of another algorithm (as shown in the previous sections), this doesn't mean unfortunately that it is possible to construct bounds that are "better" too.

We will illustrate this phenomenon with some examples (see section 6), However, we want to give without proof some general statements about upper and lower bounds first.

Lemma 2. For  $U_{\beta}$ ,  $\beta \in T$ , the sequence

$$\mathbf{x}_{n}^{\beta} := \mathbf{x}_{n-1}^{\beta} + \frac{1}{1 - c(\beta)} \|\mathbf{x}_{n}^{\beta} - \mathbf{x}_{n-1}^{\beta}\|_{\infty}, \quad n \in \mathbb{N}$$

yields monotone nonincreasing upper bounds for  $u_{\mbox{f}_{\star}}$  . Where  $c(\beta)$  is the contraction factor corresponding with  $U_{\beta}$  and where

$$\|\mathbf{x} - \mathbf{y}\|_{\infty} := \max_{i} |\mathbf{x}(i) - \mathbf{y}(i)|$$
,  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{N}$ .

Furthermore,

$$\lim_{n\to\infty} \bar{x}_n^{\beta} = u_f.$$

Lemma 3. For  $U_{\beta}^{(\lambda)}$ ,  $\beta \in T$ ,  $\lambda \in \mathbb{N}$ , the sequence

$$\bar{\mathbf{x}}_{\mathbf{n}}^{\beta\lambda} := \min \left\{ \bar{\mathbf{x}}_{\mathbf{n}-1}^{\beta\lambda} \right., \, \mathbf{x}_{\mathbf{n}-1}^{\beta\lambda} + \frac{1}{1-c(\beta)} \left\| \mathbf{U}_{\beta} \mathbf{x}_{\mathbf{n}-1}^{\beta\lambda} - \mathbf{x}_{\mathbf{n}-1}^{\beta\lambda} \right\|_{\infty}$$

yields monotone nonincreasing upper bounds for  $\mathbf{u}_{\mathbf{f}_{\star}}$ . Furthermore

$$\lim_{n\to\infty} \bar{x}_n^{\beta\lambda} = u_f.$$

It is also possible to construct a monotone nondecreasing sequence of lower bounds for u and u  $f_n^{\beta}$ , and so for u . Such a sequence can be formed trivially by using the  $x_n^{\beta}$ ,  $x_n^{\beta\lambda}$ ,  $\beta\in T$ , respectively.

We will now give sequences of lower bounds that might be used for the several methods described in the previous sections.

<u>Lemma 4</u>. For  $U_{\beta}$ ,  $\beta \in T$ , the sequence

$$\underline{\mathbf{x}}_{n}^{\beta} := \mathbf{x}_{n-1}^{\beta} + \frac{1}{1 - \delta(\beta)} \|\mathbf{x}_{n}^{\beta} - \mathbf{x}_{n-1}^{\beta}\|_{-\infty}$$

yields monotone nondecreasing lower bounds for u  $_{\mbox{\scriptsize f}_{n}^{\beta}}$  and so for u  $_{\mbox{\scriptsize f}_{n}^{\ast}}$  and

$$\lim_{n\to\infty} \underline{x}_n^{\beta} = u_{f_{\star}},$$

where  $\delta(b) = \alpha$ ;  $\delta(h) = \alpha^N$ ,  $\delta(0) = \gamma(\omega^+(f_n^0))$ ,  $\delta(s) = \gamma(\omega^*)$ ,  $\delta(h0) = \gamma(\omega^+(f_n^{h0}))^N$ ,  $\delta(hs) = (\gamma(\omega^*))^N$  and

$$\|\mathbf{x}_{n}^{\beta} - \mathbf{x}_{n-1}^{\beta}\|_{-\infty} := \min_{i \in S} (\mathbf{x}_{n}^{\beta}(i) - \mathbf{x}_{n-1}^{\beta}(i))$$
.

Lemma 5. For  $U_{\beta}^{(\lambda)}$ ,  $\beta \in T$ ,  $\lambda \in \mathbb{N}$ , the sequence

$$\underline{\mathbf{x}}_{\mathbf{n}}^{\beta\lambda} := \mathbf{x}_{\mathbf{n}-1}^{\beta\lambda} + \frac{1}{1-\delta(\beta)} \|\mathbf{U}_{\beta}\mathbf{x}_{\mathbf{n}-1}^{\beta\lambda} - \mathbf{x}_{\mathbf{n}-1}^{\beta\lambda}\|_{-\infty}$$

yields monotone nondecreasing lower bounds for u  $_{f_{n}^{\beta\lambda}}$  and so for u  $_{f_{\star}}$  ; furthermore

$$\lim_{n\to\infty} \underline{x}_n^{\beta\lambda} = u_{f_{\star}}.$$

For all the bounds we have a monotone convergence to  $u_{f_{\star}}$ . So each number of the indicated set of algorithms can be used to estimate the optimal policy  $f_{\star}$  and the corresponding value vector  $u_{f_{\star}}$ .

The examples in section 6 show that a choice for a specific algorithm may depend on the problem under consideration.

#### § 6. Examples.

In this section we will give two simple examples to illustrate that the decision which algorithm has to be chosen, might depend on the problem under consideration.

Example 1. In this example we compare the distance between the upper and lower bounds in the n-th iteration step of algorithm I and this distance in the n-th iteration step of algorithm IV.

Consider a two state problem with in each state only one possible decision. Let the matrix of transition probabilities be given by:

$$P := \begin{pmatrix} p & 1-p \\ 1-p & p \end{pmatrix}$$

and the reward vector r by: r :=  $\begin{bmatrix} r_1 \\ r_2 \end{bmatrix}$ ,  $(r_1 \ge r_2)$ , with discount factor  $\alpha$ . Then algorithm I yields

$$\mathbf{x}_n^b = \sum_{k=0}^{n-1} \alpha^k \mathbf{P}^k \mathbf{r} ,$$

so

$$x_{n}^{b} - x_{n-1}^{b} = \alpha^{n-1} p^{n-1} r = \alpha^{n-1} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} r + (-\alpha(1-2p))^{n-1} \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix} r.$$

This yields:

$$D^{b}(n) := \|x_{n}^{b} - x_{n-1}^{b}\|_{\infty} - \|x_{n}^{b} - x_{n-1}^{b}\|_{-\infty} = (\alpha | 1 - 2p |)^{n-1}(r_{1} - r_{2}).$$

Using algorithm IV yields in a similar way

$$\mathbf{D}^{\mathbf{S}}(\mathbf{n}) := \|\mathbf{x}_{\mathbf{n}}^{\mathbf{S}} - \mathbf{x}_{\mathbf{n}-1}^{\mathbf{S}}\|_{\infty} - \|\mathbf{x}_{\mathbf{n}}^{\mathbf{S}} - \mathbf{x}_{\mathbf{n}-1}^{\mathbf{S}}\|_{-\infty} = \left(\frac{\alpha(1-p)}{1-\alpha p}\right)^{n-1} \frac{1}{1-\alpha p} (\mathbf{r}_1 - \mathbf{r}_2) .$$

Let An be defined by:

$$A_{n} := \frac{D^{b}(n)}{D^{s}(n)} = c \frac{|1 - 2p|^{n-1} (1 - \alpha p)^{n-1}}{(1 - p)^{n-1}},$$

where c is a constant which is independent of n.

Then

$$\lim_{n\to\infty} Q_n = \infty \quad \text{if} \quad p > p_1 = \left(\frac{3+\alpha}{4\alpha}\right) - \sqrt{\left(\frac{3+\alpha}{4\alpha}\right)^2 - 1}$$

$$\lim_{n\to\infty} Q_n = 0 \quad \text{if} \quad p < p_1.$$

For this problem this leads to the conclusion that algorithm I is preferable if  $p < p_1$ .

Example 2. Consider a two state problem with  $K(1) = \{1\}$ ,  $K(2) = \{1,2\}$ , a = 0,9, and

$$p_{11}^{1} = 1$$
,  $p_{12}^{1} = 0$ ,  $r^{1}(1) = 2$   
 $p_{21}^{1} = 0$ ,  $p_{22}^{1} = 1$ ,  $r^{1}(2) = 2$   
 $p_{22}^{2} = 1$ ,  $p_{22}^{2} = 0$ ,  $r^{2}(2) = 1,9$ .

The Hastings algorithm II will start in state 2 with the suboptimal decision 2, while (MacQueen) algorithm I starts with selecting the optimal decision 1. Furthermore, the upper and lower bounds corresponding to algorithm I are equal, which means that the optimal values  $\mathbf{u}_{\mathbf{f}}$  are known in one step.

#### References.

- [1] Blackwell, D., Discounted dynamic programming. Am. Math. Stat. 36 (1965), 226-234.
- [2] Denardo, E.V., Contraction mappings in the theory underlying dynamic programming. SIAM Review 9 (1967), 165-177.
- [3] Hastings, N., Some notes on Dynamic Programming and Replacement.
  Opl. Res. Q. 19 (1968), 453-464.
- [4] Howard, R., Dynamic Programming and Markov Processes. M.I.T. Press, Cambridge (1960).
- [5] MacQueen, J., A modified dynamic programming method for Markovian decision problems. J. Math. An. Appl. 14 (1966), 38-43.
- [6] MacQueen, J., A test for suboptimal actions in Markovian decision problems. O.R. 15 (1967), 559-561.
- [7] Nunen, J. van, A set of successive approximation methods for discounted Markovian decision problems. Memorandum COSOR 73-09, Department of Math., Techn. Univ. Eindhoven, Netherlands.
- [8] Reetz, D., Solution of a Markovian Decision Problem by Successive Overrelaxation. Zeitschrift Operat. Res. 17 (1973), 29-32.
- [9] Schellhaas, N., Zur Extrapolation in Markoffschen Entscheidungsmodellen mit Diskontierung (1973), preprint nr. 84. Technische Hochschule Darmstadt, West-Germany.
- [10] Wessels, J. and Nunen, J. van, Discounted semi Markov decision processes: linear programming and policy iteration. Memorandum COSOR 74-01, Department of Math., Techn. Univ. Eindhoven, Netherlands.