

Memo.No. 783

An Introduction to
Dynamic Programming

by

Dr. Nadia Makary

June 1967

THE
LIBRARY OF THE
MUSEUM OF NATURAL HISTORY

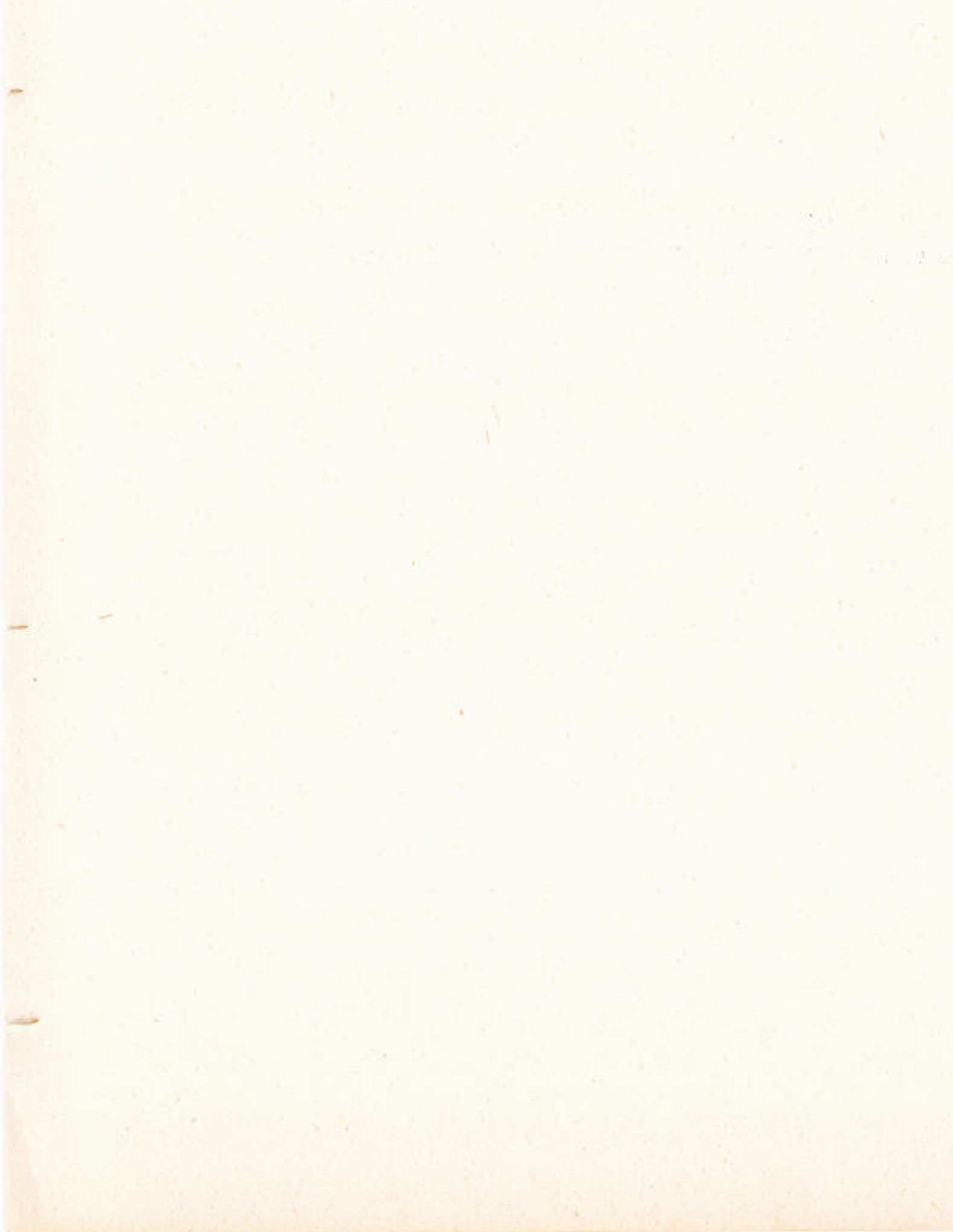
OF THE
CITY OF NEW YORK
AND
THE
LIBRARY OF THE
AMERICAN MUSEUM OF NATURAL HISTORY

18

THE
LIBRARY OF THE
MUSEUM OF NATURAL HISTORY

OF THE
CITY OF NEW YORK

"Opinions Expressed and Positions Taken by
Authors are Entirely their Own and do not Necess-
arily Reflect the Views of the Institute of Natio-
nal Planning."



- Example 2. " Gold-mining : finite-stage case ",
- Example 3. " Gold-mining : infinite stage case "
- Example 4: " The Toymaker Example "
- Example 5: " Seasonal Employment ",
- Example 6: " Inventory Problem : finite-stage ca
- Example 7: " Inventory Problem : infinite-stage c

Preface:

Dynamic programming is an approach for a multistage decision processes and finding out the optimal policy. This note is a simple introduction to this approach. It first gives a general description of situations where dynamic programming may be applied. A number of examples is given to illustrate the technique and to classify the dynamic programming techniques.

policy that maximizes (or minimizes) a predefined cost) function.

Historically, it was developed mainly through papers (1950's) as a result of his trials to solve a class of "programming" problems involving time as a significant factor. However, the dynamic programming approach is used for many static processes that can be formulated as dynamic processes.

In contrast to linear programming, there is no standard set up for dynamic programming problems. Yet, there are certain features common to all problems that can be solved by the dynamic programming approach.

A general description of the situation where the dynamic programming approach can be applied may be presented as follows:

A system may be found in one of a possible number of states. At each state there is a number of possible actions. By making any of these actions, i.e., by making a decision, the system moves from one state to another in either a deterministic or a probabilistic way. Consequently, a certain income (or cost) is earned at each state. The process continues for either finite or infinite time.

future value.

The elements of the dynamic programming problem:

S : set of states.

A : set of possible actions,

Noticethat the action may depend on the state.

q : "the law of motion" of the system. It associates each pair (s, a) a probability distribution on S . In the deterministic case $q(s'/(s, a))$ equals 1 for the specified state $s' = s$ and equals zero for any $s' \neq s$.

$i(s, a)$: the immediate return function. It determines (or cost) if the system is in state s and act a .

β : $0 \leq \beta \leq 1$; the discount factor.

The following definitions will be used:

To make a decision: is to choose one of the possible actions.

A policy : is a sequence of decisions.

An optimal policy: is a policy that maximizes (or minimizes) total expected discounted income (or cost).

Thus the dynamic programming problem as defined is solved if the structure of the optimal policy is known.

an optimal policy for the remaining stages is independent of the policy adopted in previous stages. The direct result of this principle is the development of the functional equation which gives the recurrence relation between the values of the return function in successive stages, and thus enables one to find the optimal policy for each state with n stages remaining, by first finding the optimal policy for each state with $(n-1)$ stages remaining.

The dynamic programming approach has been successfully applied to a wide variety of problems in different fields. In the following, a number of examples will be discussed. Some are given mainly to clarify the dynamic programming formulation and technique, others are presented to give an idea of some applications.

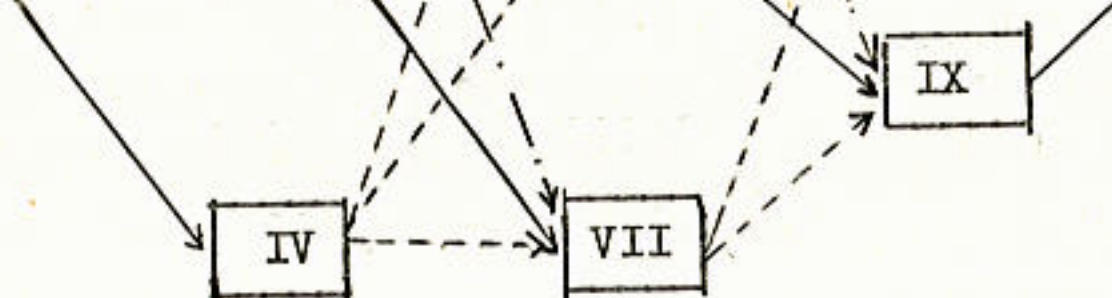
(*)

Example 1 :

(finite-stage deterministic process.)

A person wants to travel from city I to city J as quickly as possible. Owing to the long distance between the two cities, he has to make several stops before reaching his destination. At each stop he can choose the route to the next stop as shown in the diagram :

(*) This example is known in the literature by: "Stage-wise decision making".



The number of hours necessary to go from one next depends on the route he chooses as given in the

to from	II	III	IV
I	2	4	3

to from	V	VI	VII
II	7	4	6
III	3	2	4
IV	4	1	5

to from	VI
V	
VI	
VII	

Which routes should he choose in order to go in the minimum number of hours?

Trail and ~~error~~ may be used for solving this the dynamic programming approach provides an easier and solution.

For example: $a(II) = a(III) = a(IV) = V$ or VI
 $i(s, a)$ = the immediate "cost" of being in state s and
action $a(s)$ (i.e., the member of hours needed
to $a(s)$) .

For example : $i(V, IV) = 4$.

Take $\beta = 1$.

Notice that the traveller has to make four
decisions. Each time he is confronted by a decision, v
called a "stage". So this problem is a 4-stage decision

The stages will be numbered in a backward order.
at the first stage the person is either in state VIII or
IX and he should make the decision X, while at the fourth
he is in state I and has to choose one of the actions
IV.

Now, what is the optimal policy? In other words
is the sequence of routes that minimizes the total "cost".

The solution :

Let $f_n(s, a(s))$, $n = 1, 2, 3, 4$ denote the
of being in state s (at the n^{th} stage), taking action a
following the optimal policy in the remaining $(n-1)$ stages.

$$f_1(VIII) = 3,$$

$$f_1(IX) = 4,$$

$$\text{with } a^*(VIII) = a^*(IX) = X$$

Suppose now that the traveller is in state VI decides to go to VIII his total "cost" will be :

$$\begin{aligned} f_2(VI, VIII) &= i(VI, VIII) + f_1(VIII) \\ &= 6 + 3 = 9. \end{aligned}$$

But if he chooses to go to IX, his total cost

$$\begin{aligned} f_2(VI, IX) &= i(VI, IX) + f_1(IX) \\ &= 3 + 4 = 7. \end{aligned}$$

$$\text{Thus, } f_2(VI) = 7 \text{ and } a^*(VI) = IX.$$

Similarly, $f_2(s)$ and $a^*(s)$ can be calculated for each state in the second stage. The consequent results are:

		$f_2(s, a(s))$			
S	$a(s)$	VIII	IX	$f_2(s)$	$a^*(s)$
V		4	8	4	VIII
VI		9	7	7	IX
VII		6	7	6	VIII

With this information about the optimal policy for each state in the second stage, the optimal decision in each state in the third stage can be found. For example, in state I

		$f_3(s, a(s))$			$f_3(s)$	$a^*(s)$
$s \backslash a(s)$		V	VI	VII		
II		11	11	12	11	V or V
III		7	9	10	7	V
IV		8	8	11	8	V or V

In exactly the same way, the optimal decision state in the fourth stage is found; from the following be either III or IV:

		$f_4(s, (s))$			$f_4(s)$	$a^*(s)$
$s \backslash a(s)$		II	III	IV		
I		13	11	11	11	III or IV

These results show that at the initial state should go to either III or IV. If he chooses III then from there to V. From V he should go to VIII and from there to X. His total "cost" (total number of hours he spends to get from I to X) in this policy is 11 and it is the minimum possible cost. There are two more optimal routes that he can follow and still save 11 hours. These alternative routes are:

and

$I \longrightarrow IV \longrightarrow V \longrightarrow VIII$
 $I \longrightarrow IV \longrightarrow VI \longrightarrow IX$

the gold mine and remain in working order, and a probability P_1 that it will mine no gold and be damaged beyond repair, depending probabilities if it is used to mine gold in G and $(1-P_2)$ with a fraction r_2 ($0 < P_2, r_2 < 1$).

The process begins by using the machine in either F or G. If the machine is undamaged, another choice for using the machine in either F or G is made. The process continues in this manner a finite number of times if the machine is undamaged, otherwise the process ends when the machine is damaged.

What sequence of choices maximizes the amount of gold mined before the end of the process?

The dynamic programming formulation:

$$S = \left\{ s = (\alpha, \gamma) : \alpha = (1-r_1)^k x, \gamma = (1-r_2)^\ell y ; \right. \\ \left. k, \ell = 0, \dots, n ; k + \ell = n ; n = 0, 1, \dots, N \right\}$$

$$A = \{ F, G \}, \text{ where } : a = F \text{ means that mine F is to be mined,} \\ : a = G \text{ means that mine G should be mined.}$$

$$q : \begin{aligned} q(s' / (\alpha, \gamma), F) &= \begin{cases} P_1 & \text{if } s' = ((1-r_1)\alpha, \gamma) \\ 1-P_1 & \text{if } s' = (\alpha, \gamma) \end{cases} \\ q(s' / (\alpha, \gamma), G) &= \begin{cases} P_2 & \text{if } s' = (\alpha, (1-r_2)\gamma) \\ 1-P_2 & \text{if } s' = (\alpha, \gamma) \end{cases} \end{aligned}$$