

EXISTENCE OF OPTIMAL POLICIES IN STOCHASTIC DYNAMIC PROGRAMMING

BY

LAWRENCE D. BROWN AND BHARAT T. DOSHI
(NEW BRUNSWICK, NEW JERSEY)

Abstract. This paper* deals with a general discrete-time stochastic dynamic programming model. Under rather general conditions on the cost functions and the law of motion it is shown that there exists a fully optimal Borel measurable policy, that is, a policy which is optimal for future at every stage and every possible history of the process up to that stage. For the stationary dynamic programming model this implies the existence of a fully optimal stationary policy.

1. INTRODUCTION

Stochastic dynamic programming is a problem of sequential decisions under uncertainty. The basic elements of such a problem can be described as follows: We have an observable stochastic process $\{x_t: t \in T = \{1, 2, \dots\}\}$, $x_t \in X_t$, $t \in T$. At each stage t an action a_t can be chosen from the action space A_t . This selection can be based only on the history $(x_1, a_1, x_2, \dots, a_{t-1}, x_t)$ of the process $\{x_t, a_t: t \in T\}$ up to time t . Given the history $(x_1, a_1, x_2, \dots, a_{t-1}, x_t)$ and the action a_t at time t , the state x_{t+1} is determined by the probability measure

$$q_t(\cdot | x_1, x_2, \dots, x_t, a_1, a_2, \dots, a_t).$$

* The first-named author, from Cornell University, was supported in part by NSF Grant MCS 75023343-A01. The second-named author, from Bell Laboratories, was supported in part by Summer Research Fellowship, Rutgers University.

This work was done while the authors were at Rutgers University.

Also a cost $l_t(x_1, x_2, \dots, x_t, a_1, a_2, \dots, a_t)$ is incurred at time t . For a given realization $\{x_t, a_t: t \in T\}$ the loss is

$$\sum_{t \in T} l_t(x_1, x_2, \dots, x_t, a_1, a_2, \dots, a_t).$$

A *policy* δ is a measurable rule of selecting actions. For a given initial state x_1 and a policy δ the *risk* $R(x_1, \delta)$ is defined to be the expectation of the loss function under δ . For a given initial measure μ and a policy δ , the risk is the expectation of $R(x_1, \delta)$ with respect to μ . An x_1 -*optimal policy* is the one minimizing $R(x_1, \delta)$ over all policies. A μ -*optimal policy* is a policy which minimizes $R(\mu, \delta)$ over all policies. An *optimal policy* minimizes $R(x_1, \delta)$ over all policies for all x_1 . Finally, a policy that is optimal at each stage of the process will be called *fully optimal*.

Under appropriate assumptions we show that there exists a fully optimal policy. This of course implies that for each initial measure μ there exists a μ -optimal policy, a result which was previously proved by Schäl [13] under assumptions similar to ours.

When the state and action spaces are independent of time and the one-step cost function l_t is independent of t up to a multiple α^t , the resulting problem is called *stationary dynamic programming*. Stationary dynamic programming was extensively studied by Blackwell [2], [3] and Strauch [16]. One of their main concerns was to establish conditions under which a stationary optimal policy exists. Both gave examples in which a Borel measurable stationary optimal policy does not exist. Under our assumptions on the action space and the one-step cost function this difficulty does not arise and we can show the existence of a stationary fully optimal policy.

For the stationary dynamic programming Freedman [7], Furukawa [8] and Schäl [14] also give sufficient conditions for the existence of a Borel measurable stationary optimal policy. Freedman [7] (condition (3)) and Schäl [14] (condition (W)) assume continuity of q and lower semi-continuity of l in both the state and action variables. Furukawa [8] and Schäl [14], (S), assume continuity of q and lower semi-continuity of l in the action only. Their conditions on l are similar, but less explicit than ours. Thus, in the special case of stationary dynamic programming, Schäl's Theorem 15.2 is similar to our Theorem 5.1. (It is better in the sense that his regularity condition (S2) or (S2)' is weaker than our corresponding Assumption 2.12 when applied to the stationary case.) Under the assumptions of Freedman [7] and Schäl [14], (W), it is possible to convert the general non-stationary model into a stationary one by considering the entire history up to any stage as the state variable at that stage. Thus their results give the existence of a fully optimal Borel measurable policy in the non-stationary case under the appropriate regularity conditions on l and q . However, these conditions do

not hold in our model. Such a conversion is not possible in (S) of [14] or in [8] because conditions (S2) and (S2)' of Schäl [14] will be violated in the resulting stationary model. In our model, Assumption 2.12 would not hold if we were to attempt such a conversion. This necessitates an explicit treatment of a non-stationary model. Such an explicit treatment of a non-stationary model may be found in Hinderer [9], Kertz and Nachman [10] and Schäl [15] among which only Kertz and Nachman [10] prove the existence of a fully optimal policy. However, their assumptions do not imply, nor are implied by ours.

Our results and proofs rely heavily on a general decision problem studied by Brown [4] and the measurable selections of extrema investigated by Brown and Purves [5].

2. CHARACTERISTICS AND FORMULATION

We now characterize a (stochastic) dynamic programming problem and specify assumptions on its various components.

Let $T_k = \{k, k+1, \dots\}$ and $T = T_1$.

2.1. State space. For each $t \in T$, the *state space* X_t at stage t is assumed to be a Borel subset of a complete separable metric space. Let X_t be endowed with the usual topology and let \mathcal{B}_t denote the Borel σ -algebra on X_t .

2.2. Sample space. The *sample space* is

$$X = \prod_{t \in T} X_t.$$

Let

$$X_{(t)} = \prod_{s \leq t} X_s, \quad \mathcal{B}_{(t)} = \prod_{s \leq t} \mathcal{B}_s \quad \text{and} \quad \mathcal{B} = \prod_{t \in T} \mathcal{B}_t = \mathcal{B}_{(\infty)}.$$

2.3. Action space. The set of available actions at stage t is an *action space* A_t , which is assumed to be a compact subset of a complete separable metric space. Let τ_t and \mathcal{A}_t denote the usual topology and the Borel σ -algebra on A_t , respectively. Let

$$A = \prod_{t \in T} A_t$$

with the product σ -algebra \mathcal{A} . For $t \in T$, let

$$A_{(t)} = \prod_{s \leq t} A_s \quad \text{and} \quad \mathcal{A}_{(t)} = \prod_{s \leq t} \mathcal{A}_s.$$

We will use π_t and $\pi_{(t)}$ to denote the natural projection maps on X , A , etc.

2.4. Law of motion. The *law of motion* is a sequence $q = \{q_t\}$ such that q_t is a measurable conditional distribution (Markov kernel) on \mathcal{B}_{t+1} given $\mathcal{B}_{(t)} \times \mathcal{A}_{(t)}$. Formally,

(i) $q_t(\cdot | x, a)$ is a version of the conditional probability distribution on \mathcal{B}_{t+1} given $\mathcal{B}_{(t)} \times \mathcal{A}_{(t)}$.

(ii) $q_t(B | \cdot, \cdot)$ is $(\mathcal{B}_{(t)} \times \mathcal{A}_{(t)})$ -measurable for all $B \in \mathcal{B}_{t+1}$.

The interpretation is that, given $x_{(t)} = \pi_{(t)}(x)$ and $a_{(t)} = \pi_{(t)}(a)$, x_{t+1} is selected according to the measure $q_t(\cdot | x, a)$. (Note that $q_t(\cdot | x, a)$ depends on x and a only through $x_{(t)}$ and $a_{(t)}$. We will adopt similar notational conventions later in the paper.)

For any probability measure μ on the initial state x_1 , the law of motion q will then define a conditional distribution F_μ on (X, \mathcal{B}) given \mathcal{A} according to the natural formula

$$\begin{aligned} F_\mu \left(\left(\prod_{t=1}^n B_t \right) \times \left(\prod_{t=n+1}^{\infty} X_t \right) \middle| a \right) \\ = \int_{B_1} \dots \int_{B_{n-1}} \int_{B_n} q_{n-1}(dx_n | x, a) q_{n-2}(dx_{n-1} | x, a) \dots q_1(dx_2 | x, a) \mu(dx_1). \end{aligned}$$

Let μ_{x_1} denote the distribution degenerate at $x_1 \in X_1$. We will write $F_{x_1} = F_{\mu_{x_1}}$.

2.5. Policies. A (*sequential*) *policy* is a measurable non-anticipative rule of selecting actions. Formally, $\delta = \{\delta_t\}$ is a *policy* if

- (i) $\delta_t(\cdot | x, a)$ is a probability measure on \mathcal{A}_t for all $x \in X$ and $a \in A$, $t \geq 2$;
- (ii) $\delta_t(C | \cdot, \cdot)$ is $(\mathcal{B}_{(t)} \times \mathcal{A}_{(t-1)})$ -measurable for each $C \in \mathcal{A}_t$, $t \geq 2$;
- (iii) $\delta_1(\cdot | x)$ is a probability measure on \mathcal{A}_1 for all $x \in X$;
- (iv) $\delta_1(C | \cdot)$ is \mathcal{B}_1 -measurable for all $C \in \mathcal{A}_1$.

The interpretation is that, if $x_1, a_1, \dots, a_{t-1}, x_t$ is the observed history up to stage t , then a_t is chosen according to the measure $\delta_t(\cdot | x, a)$. Let $\mathcal{D} = \{\delta\}$ denote the set of all policies.

2.6. Observed process. The preceding suffices to guarantee the existence of a well-defined stochastic process of the observed states and action when the initial measure is μ and a policy δ is used. This process is defined on $X \times A$ with the corresponding σ -algebra. The probability of a cylinder set

$$(B_1, C_1) \times (B_2, C_2) \times \dots \times (B_n, C_n) \times (X_{n+1}, A_{n+1}) \times (X_{n+2}, A_{n+2}) \dots$$

is given by

$$\begin{aligned} A_{\mu, \delta} \left(\prod_{t=1}^n (B_t, C_t) \times \prod_{t>n} (X_t, A_t) \right) \\ = \int_{B_1} \int_{C_1} \dots \int_{B_n} \int_{C_n} \delta_n(da_n | x, a) q_{n-1}(dx_n | x, a) \dots \delta_1(da_1 | x) \mu(dx_1). \end{aligned}$$

Similarly, if μ is a probability measure on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$, $k \in T$, then define $\Delta_{\mu, \delta}^k$ by

$$\begin{aligned} & \Delta_{\mu, \delta}^k \left(\prod_{t=1}^n (B_t, C_t) \times \prod_{t>n} (X_t, A_t) \right) \\ &= \int_{B_1} \int_{C_1} \dots \int_{B_n} \int_{C_n} \delta_n(da_n | x, a) q_{n-1}(dx_n | x, a) \dots \delta_k(da_k | x, a) \mu(d(x_{(k)}, a_{(k-1)})). \end{aligned}$$

Note that $\Delta_{\mu, \delta}^1 = \Delta_{\mu, \delta}$ and also that $\Delta_{\mu, \delta}^k$ depends on δ only through $\{\delta_t: t \geq k\}$.

2.7. Costs. The cost incurred at stage t is $l_t(x_{(t)}, a_{(t)})$ if the observed history is $(x_1, a_1, \dots, a_{t-1}, x_t)$ and a_t is chosen at stage t .

2.8. Assumption. For $t \in T$

- (i) $l_t(\cdot, \cdot)$ is $(\mathcal{B}_{(t)} \times \mathcal{A}_{(t)})$ -measurable;
- (ii) $l_t(\cdot, \cdot)$ is non-negative extended real valued;
- (iii) $l_t(x_{(t)}, \cdot)$ is lower semi-continuous for each $x_{(t)} \in X_{(t)}$.

In most of what follows the non-negativity of $l_t(\cdot, \cdot)$ can be replaced by the weaker assumption that $l_t(\cdot, \cdot)$ is bounded from below or by the even weaker assumption that $\{l_t(\cdot, \cdot)\}$ is uniformly integrable.

2.9. Loss function. The loss function $L: X \times A \rightarrow R^+$ is defined by

$$L(x, a) = \sum_{t=1}^{\infty} l_t(x_{(t)}, a_{(t)}) \quad (x \in X, a \in A).$$

Also define $L_k^n: X \times A \rightarrow R^+$ by

$$L_k^n(x, a) = \sum_{t=k}^n l_t(x_{(t)}, a_{(t)}) \quad (1 \leq k \leq n \leq \infty, k < \infty).$$

Let $L_k = L_k^\infty$ and $L^n = L_1^n$ ($1 \leq k \leq n \leq \infty$).

2.10. Risk function. Let $1 \leq k \leq n \leq \infty$ as above. For a measure μ on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$ and a policy $\delta \in \mathcal{D}$, the conditional risk R_k^n is defined to be the expectation of the loss function $L_k^n(\cdot, \cdot)$ computed under the measure $\Delta_{\mu, \delta}^k$:

$$R_k^n(\mu, \delta) = E_{\Delta_{\mu, \delta}^k} (L_k^n(x, a)) = \int_{X \times A} L_k^n(x, a) \Delta_{\mu, \delta}^k(dx, da).$$

Let $R_k = R_k^\infty$, $R^n = R_1^n$ and $R = R_1$. Then $R(\mu, \cdot)$ is called the risk for the initial measure μ (on \mathcal{B}_1). In particular, if μ_{x_1} is the probability measure degenerate at $x_1 \in X_1$, then let $R(x_1, \delta) = R(\mu_{x_1}, \delta)$.

Given μ on \mathcal{B}_1 , define $\Delta_{\mu, \delta|k}$ to be the projection of $\Delta_{\mu, \delta}$ on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$. It will be of special interest to consider $R_k(\Delta_{\mu, \delta|k}, \delta)$, for this can be interpreted as the conditional risk from stage k (inclusive) onward, given that the initial measure was μ , the decision rules $\delta_1, \delta_2, \dots, \delta_{k-1}$ have already been used, and that the rules $\delta_k, \delta_{k+1}, \dots$ will be used from stage k onward.

2.11. Optimal policies. Let $k \in T$ and let μ be a probability measure on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$. If

$$R_k(\mu, \delta') = \inf_{\delta \in \mathcal{D}} R_k(\mu, \delta),$$

then δ' is called μ -optimal at stage k . If $\mu = \mu_{(x_{(k)}, a_{(k-1)})}$, then δ' is called $(x_{(k)}, a_{(k-1)})$ -optimal at stage k . A policy δ is called optimal at stage k if it is optimal for $(x_{(k)}, a_{(k-1)})$ at stage k for all $(x_{(k)}, a_{(k-1)}) \in X_{(k)} \times A_{(k-1)}$. A policy is called fully optimal if it is optimal at stage k for all $k \in T$.

For $(x_{(k)}, a_{(k-1)}) \in X_{(k)} \times A_{(k-1)}$ let

$$R_k(x_{(k)}, a_{(k-1)}) = \inf_{\delta \in \mathcal{D}} R_k(\mu_{(x_{(k)}, a_{(k-1)})}, \delta).$$

Define $R_k^n(\cdot; \cdot)$ similarly relative to the loss function L_k^n .

2.12. Assumption. (a) The law of motion $q = \{q_t\}$ is dominated. That is, there exists a law of motion $v = \{v_t\}$ which is independent of $a \in A$ and such that $q_t \ll v_t$. Let

$$\varphi_t(\cdot | x, a) = \frac{dq_t(\cdot | x, a)}{dv_t(\cdot | x)} \quad (a \in A, t \in T).$$

We also assume that $\varphi_t(\cdot | \cdot, \cdot)$ is $(\mathcal{B}_{t+1} \times \mathcal{B}_{(t)} \times \mathcal{A}_{(t)})$ -measurable.

(b) Assume also that for each $t \in T$, $x \in X$, the map

$$A \rightarrow L_1(X_{t+1}, \mathcal{B}_{t+1}, v_t(\cdot | x))$$

defined by $a \rightarrow \varphi_t(\cdot | x, a)$ is continuous. Note that this map depends on a only through $a_{(t)} = \pi_{(t)}(a)$.

Let μ denote a probability measure on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$ and let $\hat{\mu}$ denote its projection (marginal measure) on $\mathcal{B}_{(k)}$. Let V_μ^k denote the measure on \mathcal{B} generated by $\hat{\mu}$ through stage k followed by v . Let $\varphi_\mu^k(\cdot | a)$ denote the conditional measure on \mathcal{B} defined by

$$\begin{aligned} & \varphi_\mu^k(\mathcal{B}_{(k)} \times \prod_{t=k+1}^n \mathcal{B}_t \times \prod_{t=n+1}^\infty X_t | a) \\ &= \int_{A_{(k-1)}} \int_{\mathcal{B}_{(k)}} \int_{\mathcal{B}_{k+1}} \dots \int_{\mathcal{B}_n} q_{n-1}(dx_n | x, a) \dots q_k(dx_{k+1} | x, a) \mu(d(x_{(k)}, a_{(k-1)})). \end{aligned}$$

Note that $\varphi_\mu^k(\cdot | a)$ depends on a only through a_k, a_{k+1}, \dots . Assumption 2.12 (a) guarantees that $\varphi_\mu^k(\cdot | a) \ll V_\mu^k(\cdot)$ for all $a \in A$, $k \in T$, μ on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$. (It also guarantees that the projection of $\Delta_{\mu, \delta}^k$ on \mathcal{B} is dominated by V_μ^k .) Let

$$f_\mu^k(\cdot | a) = \frac{d\varphi_\mu^k(\cdot | a)}{dV_\mu^k(\cdot)}.$$

Then Assumption 2.12 (b) implies that for each $k \in T$ the map

$$N_k(\mu): A \rightarrow L_1(X, \mathcal{B}, V_\mu^k)$$

defined by $a \rightarrow f_\mu^k(\cdot | a)$ is continuous for each $k \in T$ and μ on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$.

Note that $\varphi_\mu^1 = F_\mu$, and so, in particular, the map

$$N_1(\mu): a \rightarrow \frac{dF_\mu}{dV_\mu^1} = f_\mu^1$$

is continuous for each initial measure μ on \mathcal{B}_1 .

The assumption above is a significant restriction on the family g . However, the necessity of some assumption similar to 2.12 (b) is demonstrated in Example 3.4.

2.13. Topologizing the policies. Corresponding to each policy $\delta = \{\delta_t\} \in \mathcal{D}$ there is a conditional measure on A , given \mathcal{B} , defined by

$$\begin{aligned} \delta(C_1 \times C_2 \times \dots \times C_t \times A_{t+1}, \dots | x) \\ = \int_{c_1} \int_{c_2} \dots \int_{c_{t-1}} \delta_t(C_t | x, a) \delta_{t-1}(da_{t-1} | x, a) \dots \delta_1(da_1 | x). \end{aligned}$$

For each $k \in T$ and μ on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$ define the topology $\tau_k(\mu)$ on \mathcal{D} as the weak topology generated by the functions

$$\delta \rightarrow \int \delta(da | x) f(x) c(a) V_\mu^k(dx)$$

for $f \in L_1(X, \prod_{t \in T_k} \mathcal{B}_t, V_\mu^k)$ and $c \in C(\prod_{t \in T_k} A_t)$, where $C(\prod_{t \in T_k} A_t)$ denotes the set of continuous functions on $\prod_{t \in T_k} A_t$. Note that $\tau_k(\mu)$ really depends only on the projection measure $\hat{\mu}$ on $\mathcal{B}_{(k)}$. Let $\tau(\mu) = \tau_1(\mu)$ and $\tau(x_1) = \tau(\mu_{x_1})$, $x_1 \in X_1$.

These topologies are not Hausdorff. For an intensive study of \mathcal{D} with these topologies it is useful to define the space $\mathcal{D}_k(\mu)$ of equivalence classes of \mathcal{D} with the topology $\tau_k(\mu)$. This is done in [4]. It may help the reader to visualize \mathcal{D} , $\tau(\mu)$, etc., in this manner but we will not explicitly need this terminology in the sequel.

The following important results may be directly deduced from Theorems 3.6 and 3.14 of [4].

2.14. THEOREM. For each $k \in T$, μ on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$,

- (a) the spaces $\mathcal{D}_k, \tau_k(\mu)$ are compact,
- (b) the maps $R_k(\mu, \cdot): \mathcal{D} \rightarrow [0, \infty]$ are lower semi-continuous relative to $\tau_k(\mu)$.

It is implicit in Theorem 2.14 (b) that $R_k(\mu, \delta)$ depends only on the $\tau_k(\mu)$ equivalence class of $\delta \in \mathcal{D}$. For $k = 1$ this latter fact follows from the expression

$$R(\mu, \delta) = \int L(x, a) f_\mu^1(x|a) \delta(da|x) V_\mu(dx)$$

and from Theorem 3.10 in [4]. This expression is derived in [4], Proposition 2.2. Similar expressions are valid for $R_k(\mu, \delta)$ but are not explicitly needed in the sequel. (The above result is comparable to results in [13].)

3. POLICIES OPTIMAL FOR μ

3.1. THEOREM. *For each $k \in T$ and each μ on $\mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)}$ there exists an optimal policy for μ at stage k . In particular, for each $x_1 \in X_1$ there is an x_1 -optimal policy at stage 1.*

Proof. The theorem follows directly from Theorem 2.14 since any lower semi-continuous function on a compact set achieves its infimum.

This theorem leaves many questions unanswered. For example, if μ is a given initial measure on \mathcal{B}_1 and δ is μ -optimal at stage 1, then is δ an x_1 -optimal policy at stage 1 for μ -almost every x_1 ? We have been able to adapt methods developed in Eaton [6] to answer this question (and other similar ones) in the affirmative when the measures $\{V_{\mu_{x_1}}^1 : x_1 \in X_1\}$ form a dominated family. However, the optimality results proved in the next section by a different method are stronger than the results we can prove using Eaton's methods, and so we will omit the detailed argument.

The following technical result will later be useful:

3.2. LEMMA. $R(x_1) = \lim_{n \rightarrow \infty} R^n(x_1)$.

Proof. $R^n(x_1, \delta)$ is the risk in the problem with the loss function

$$L^n(x, a) = \sum_{t=1}^n l_t(x_{(t)}, a_{(t)}).$$

By Theorem 3.1 there exists an optimal policy δ_n for each such problem. Since $\mathcal{D}, \tau(x_1)$ is compact, there exists a "convergent" subnet $\delta_{n'} \rightarrow \delta$. (Actually δ is determined only up to its $\tau(x_1)$ equivalence class and $\delta_{n'} \rightarrow \delta$ only in this sense.) By Theorem 2.14 (b) we have

$$\liminf_{(n')} R^m(x_1, \delta_{n'}) \geq R^m(x_1, \delta) \quad \text{for each } m \in T.$$

Also

$$R(x_1) \geq R^n(x_1, \delta_n) \geq R^m(x_1, \delta_n) \quad (n \geq m).$$

Hence

$$R(x_1) \geq \lim_{m \rightarrow \infty} \limsup_{(n')} R^m(x_1, \delta_{n'}) \geq \lim_{m \rightarrow \infty} R^m(x_1, \delta) = R(x_1, \delta) \geq R(x_1).$$

This proves the lemma.

3.3. Remark. The assumption that A be compact may be relaxed. For example, it suffices for A to be locally compact with

$$\lim_{a \rightarrow \infty} L(x, a) = \infty$$

and for certain other technical conditions to be satisfied. Details may be deduced from Section 4 of [4] together with the above theorems. A similar remark is also valid for later theorems in this paper.

The forthcoming example demonstrates the important role of Assumption 2.12 (b) concerning the continuity of the maps $a \rightarrow \varphi_t(\cdot | x, a)$. The example demonstrates that Theorem 3.1 is false under the weaker assumption that for fixed μ the map

$$a \rightarrow \frac{dF_\mu(\cdot | a)}{dV_\mu^1(\cdot)}$$

is continuous as a map from A to $L_1(X, \mathcal{B}, V_\mu^1)$ with the weak topology.

3.4. Example. Let S denote the trivial space consisting of one point. Let

$$X_1 = S, \quad A_1 = \left\{ \frac{1}{m} : m = 1, 2, \dots \right\} \cup \{0\},$$

$$X_2 = [0, 1], \quad A_2 = \{0, 1\} = X_3, \quad A_k = S = X_{k+1}, \quad k = 3, 4, \dots,$$

and let

$$l_3(x_{(3)}, a_{(3)}) = \begin{cases} 0 & \text{if } a_2 = x_3, \\ 1 & \text{if } a_2 \neq x_3, \end{cases} \quad l_k \equiv 0 \quad \text{for } k \neq 3.$$

In standard statistical terminology, x_3 is the unknown parameter and the problem is one of testing the simple hypothesis, $x_3 = 0$, versus the simple alternative, $x_3 = 1$, based on the observation of $x_2 \in X_2$. The prior distribution of x_3 and the distribution of x_2 , given x_3 , are determined by the choice of a_1 as specified in the sequel.

For $y \in [0, 1]$, let $b_k(y)$ denote the k -th digit in the binary expansion of y . (Adopt any convention to define $b_k(y)$ uniquely when y has a terminating binary expansion.) Let λ denote the Lebesgue measure on $[0, 1]$ and let

$$\frac{dq_1(x_2 | x, a)}{d\lambda} = \begin{cases} (2/3 - a_1) \cdot 2b_{a_1^{-1}}(x_2) + 1/3 + a_1 & \text{if } a_1 > 0, \\ 1 & \text{if } a_1 = 0, \end{cases}$$

$$q_2(\{0\} | x, a) = \begin{cases} \frac{2(2/3 - a_1)}{5/3 - a_1} & \text{if } b_{a_1^{-1}}(x_2) = 1, a_1 > 0, \\ 0 & \text{if } b_{a_1^{-1}}(x_2) = 0, a_1 > 0, \\ 2/5 & \text{if } a_1 = 0. \end{cases}$$

The choice of $\{v_i\}$ with $v_1 = \lambda$ and $v_2 = v^*$, where $v^* (\{0\}) = v^* (\{1\}) = 1/2$, yields

$$V_\mu^1 = \tau \times \lambda \times v^* \times \tau \times \tau \times \tau \times \dots,$$

where τ denotes the unique probability measure on S . Then it can easily be checked that the map $a \rightarrow dF_\mu(\cdot|a)/dV_\mu^1(\cdot)$ is continuous as a map from A to $L_1(X, \mathcal{B}, V_\mu^1)$ with the *weak topology*. It is also true that the map $a \rightarrow q_1(\cdot|x, a)$ has a similar continuity property; but note that the map $a \rightarrow q_2(\cdot|x, a)$ does not.

Consider the sequence of policies $\{\delta^{(m)} : m = 1, 2, \dots\}$ defined by

$$\begin{aligned} \delta_1^{(m)}(\{m^{-1}\}|x) &= 1, \\ \delta_2^{(m)}(\{0\}|a) &= \begin{cases} 1 & \text{if } b_m(x_2) = 1, \\ 0 & \text{if } b_m(x_2) = 0. \end{cases} \end{aligned}$$

Then a simple computation yields $R(\mu, \delta^{(m)}) = 1/6 + 1/2m$. Hence $R(\mu) \leq 1/6$, but it can be readily seen that there is no policy which satisfies $R(\mu, \delta) \leq 1/6$. So there is no optimal policy.

4. EXISTENCE OF FULLY OPTIMAL POLICIES

Now we state the main theorem of the paper.

4.1. THEOREM. *Let the assumptions of Section 2 be satisfied⁽¹⁾. Then there exists a fully optimal policy.*

Proof. For $1 \leq k \leq n < \infty$ define $S_k^n(\cdot; \cdot, \cdot)$ inductively (backwards on k) by

$$\begin{aligned} S_n^n(a_n; x_{(n)}, a_{(n-1)}) &= l_n(x_{(n)}, (a_{(n-1)}, a_n)), \\ (1) \quad S_{k-1}^n(a_{k-1}; x_{(k-1)}, a_{(k-2)}) &= l_{k-1}(x_{(k-1)}, (a_{(k-2)}, a_{k-1})) + \\ &+ \int \left\{ \inf_{a_k \in A_k} S_k^n(a_k; (x_{(k-1)}, x_k), (a_{(k-2)}, a_{k-1})) \right\} \varphi_{k-1}(x_k|x, a) v_{k-1}(dx_k|x), \end{aligned}$$

where \int denotes the lower integral. Write $S_1^n(a_1; x_1)$ instead of $S_1^n(a_1; x_{(1)}, a_{(0)})$.

The following claims will now be proved for $1 \leq k \leq n < \infty$:

- (A) $S_k^n(\cdot; \cdot, \cdot)$ is $(\mathcal{A}_k \times \mathcal{B}_{(k)} \times \mathcal{A}_{(k-1)})$ -measurable.
 - (B) $S_k^n(\cdot; x_{(k)}, a_{(k-1)})$ is lower semi-continuous on A_k for all $x_{(k)}, a_{(k-1)}$.
 - (C) There exists a measurable map $h_{k,n} : X_{(k)} \times A_{(k-1)} \rightarrow A_k$ such that
- $$(2) \quad \begin{aligned} \inf_{a_k \in A_k} S_k^n(a_k; x_{(k)}, a_{(k-1)}) &= S_k^n(h_{k,n}(x_{(k)}, a_{(k-1)}); x_{(k)}, a_{(k-1)}) \\ &= R_k^n(x_{(k)}, a_{(k-1)}). \end{aligned}$$

⁽¹⁾ For the proof of this theorem the most important are the assumptions stated in Sections 2.3, 2.8, 2.12 and 3.3.

For $k = n$, claims (A) and (B) follow from Assumption 2.8. In view of the assumptions on A_i stated in Section 2.3, the existence of $h_{n,n}$ satisfying the first part of (2) follows from the selection theorem of Brown and Purves [5]. It is clear that

$$S_n^n(h_{n,n}(x_{(n)}, a_{(n-1)}); x_{(n)}, a_{(n-1)}) \leq R_n^n(x_{(n)}; a_{(n-1)}).$$

On the other hand,

$$R_n^n(\mu_{x_{(n)}, a_{(n-1)}}(\cdot), \delta) = S_n^n(h_{n,n}(x_{(n)}, a_{(n-1)}); x_{(n)}, a_{(n-1)})$$

for the policy δ with

$$\delta_n(\cdot | x_{(n)}, a_{(n-1)}) = \mu_{h_{n,n}(x_{(n)}, a_{(n-1)})}.$$

Hence the second equality of (2) is also satisfied for $k = n$.

Let $2 \leq m \leq n$. Suppose that (A), (B) and (C) are satisfied for $m \leq k \leq n$. Then

$$(3) \quad S_{m-1}^n(a_{m-1}; x_{(m-1)}, a_{(m-2)}) = I_{m-1}(x_{(m-1)}, (a_{(m-2)}, a_{m-1})) + \int R_m^n((x_{(m-1)}, x_m); (a_{(m-2)}, a_{m-1})) \varphi_{m-1}(x_m | x, a) v_{m-1}(dx_m | x).$$

By (A) and (C), $R_m^n(\cdot; \cdot)$ is $(\mathcal{B}_{(m)} \times \mathcal{A}_{(m-1)})$ -measurable. Hence S_{m-1}^n satisfies (A) by standard Fubini-type theorems (see, e.g., [12], p. 74).

Let α_i be any convergent sequence in A_{m-1} , say $\alpha_i \rightarrow \alpha$. Let $\delta^{(i)}$ be a policy satisfying

- (i) $\delta_k^{(i)}(\cdot | x, a) = \mu_{h_{k,n}(\cdot)(x_{(k)}, a_{(k-1)})}$ ($m \leq k \leq n$),
- (ii) $\delta_{m-1}^{(i)}(\{\alpha_i\} | \cdot, \cdot) \equiv 1$.

Then, by (2) and (3),

$$(4) \quad S_{m-1}^n(\alpha_i; x_{(m-1)}, a_{(m-2)}) = R_{m-1}^n(\mu_{(x_{(m-1)}, a_{(m-2)})}, \delta^{(i)}).$$

Let δ be any $\tau_{m-1}(\mu_{(x_{(m-1)}, a_{(m-2)})})$ accumulation point of the sequence $\{\delta^{(i)}\}$ (δ exists by Theorem 2.14 (a)). Note that every accumulation point of $\{\delta^{(i)}\}$ is equivalent to δ relative to this topology. By (4) and Theorem 2.14 (b),

$$\begin{aligned} S_{m-1}^n(\alpha; x_{(m-1)}, a_{(m-2)}) &= R_{m-1}^n(\mu_{(x_{(m-1)}, a_{(m-2)})}, \delta) \\ &\leq \liminf_{i \rightarrow \infty} R_{m-1}^n(\mu_{(x_{(m-1)}, a_{(m-2)})}, \delta^{(i)}) \\ &= \liminf_{i \rightarrow \infty} S_{m-1}^n(\alpha_i; x_{(m-1)}, a_{(m-2)}). \end{aligned}$$

Hence S_{m-1}^n satisfies (B). The existence of $h_{m-1,n}$ satisfying the first equality of (2) then follows from the selection theorem of Brown and Purves [5]. The second equality of (2) then follows from the first equality and the definitions of S_{m-1}^n and R_{m-1}^n by the same reasoning as in the case $k = n$. This completes the proof of (A), (B) and (C) for $1 \leq k \leq n < \infty$.

Clearly, S_k^n is non-decreasing in n for fixed k . Let

$$S_k = \lim_{n \rightarrow \infty} S_k^n.$$

Then S_k satisfies (A) and (B). So there exist measurable functions h_k satisfying the first equality in (2) for S_k . By Lemma 3.2,

$$\begin{aligned} (5) \quad R_k((x_{(k)}, a_{(k-1)})) &= \lim_{n \rightarrow \infty} R_k^n(x_{(k)}, a_{(k-1)}) \\ &= \lim_{n \rightarrow \infty} \inf_{a_k \in A_k} S_k^n(a_k; x_{(k)}, a_{(k-1)}) \\ &= \inf_{a_k \in A_k} S_k(a_k; x_{(k)}, a_{(k-1)}) \end{aligned}$$

since A_k is compact and $\{S_k^n(\cdot; x_{(k)}, a_{(k-1)}); n = 1, 2, \dots\}$ is a non-decreasing sequence of lower semi-continuous functions. Thus S_k and R_k also satisfy the second equality in (2). It follows from the above and from the definition (1) that, furthermore,

$$(6) \quad \lim_{n \rightarrow \infty} S_k^n(h_{k,n}(x_{(k)}, a_{(k-1)}); x_{(k)}, a_{(k-1)}) = S_k(h_k(x_{(k)}, a_{(k-1)}); x_{(k)}, a_{(k-1)})$$

since A_k is compact and $\{S_k^n(\cdot; x_{(k)}, a_{(k-1)}); n = 1, 2, \dots\}$ is a non-decreasing sequence of lower semi-continuous functions. So $\{S_k\}$ satisfies the recursion (1) for all $k \geq 2$, $(x_{(k-1)}, a_{(k-2)}) \in X_{(k-1)} \times A_{(k-2)}$.

Let δ^* be the policy defined by

$$\delta_k^* (\{h_k(x_{(k)}, a_{(k-1)})\} | x, a) = 1 \quad (k \in T).$$

Then, by (1),

$$R_k^n(\mu, \delta^*) = E_{\mu, \delta^*}^k \left(\sum_{t=k}^n l_t(x_{(t)}, a_{(t)}) \right) \leq S_k(h_k(x_{(k)}, a_{(k-1)}); x_{(k)}, a_{(k-1)}) \quad (1 \leq k \leq n < \infty),$$

where $\mu = \mu_{(x_{(k)}, a_{(k-1)})}$. Hence

$$R_k(\mu_{(x_{(k)}, a_{(k-1)})}, \delta^*) \leq S_k(h_k(x_{(k)}, a_{(k-1)}); x_{(k)}, a_{(k-1)}).$$

By Lemma 3.2 and formulas (5) and (6),

$$\begin{aligned} R_k(x_{(k)}; a_{(k-1)}) &= S_k(h_k(x_{(k)}, a_{(k-1)}); x_{(k)}, a_{(k-1)}) \\ &\geq R_k(\mu_{(x_{(k)}, a_{(k-1)})}, \delta^*) \geq R_k(x_{(k)}; a_{(k-1)}). \end{aligned}$$

It follows that $R_k(\mu_{(x_{(k)}, a_{(k-1)})}, \delta^*) = R_k(x_{(k)}; a_{(k-1)})$, and hence δ^* is optimal at stage k for all $k \in T$. This proves the theorem.

Note that the fully optimal policy δ^* , derived above, is non-randomized.

5. STATIONARY CASE

In this section we consider an important special case of the general problem. Assume that the state and action spaces, the law of motion and the cost structure are time invariant (stationary), that is,

$$X_t = X_1, \quad A_t = A_1,$$

$$q_t(B_{t+1} | x, a) = q_t(B_{t+1} | x_t, a_t) = q_{t-1}(B_t | x_{t-1}, a_{t-1}) \quad (t \geq 2)$$

if $B_{t+1} = B_t \in \mathcal{B}_1$, $x_t = x_{t-1}$ and $a_t = a_{t-1}$. Moreover, for some α , $0 < \alpha \leq 1$,

$$l_t(x_{(t)}, a_{(t)}) = \alpha^{t-1} l(x_t, a_t),$$

where $l(\cdot, \cdot)$ is a $(\mathcal{B}_1 \times \mathcal{A}_1)$ -measurable non-negative function which is lower semi-continuous with respect to the action variable.

When $\alpha = 1$, this corresponds to the negative dynamic programming of Strauch [16], and when $\alpha < 1$ and $l(\cdot, \cdot)$ is bounded (not necessarily non-negative), this corresponds to the discounted dynamic programming of Blackwell [2]. One of their main concerns was to determine whether a stationary optimal policy necessarily exists, that is, whether there always exists a $\delta \in \mathcal{D}$ such that δ is optimal in \mathcal{D} and

$$\delta_t(C | x, a) = \delta_t(C | x_t) = \delta_{t-1}(C | x_{t-1}) \quad (t \geq 2)$$

for all $C \in \mathcal{A}_1$ and $x_t = x_{t-1} \in X_1$.

Blackwell [2] gave a counterexample showing that this is not the case. However, in his example, depending on the interpretation, either A_1 is not compact or $l(x, \cdot)$ is not lower semi-continuous. With our assumptions, which include compactness of A_1 (but see Remark 3.3), lower semi-continuity of $l(x, \cdot)$ and assumptions on q , we can show that a fully optimal stationary policy does exist.

The desired result follows almost immediately from Theorem 4.1. (It could also be proved independently of Theorem 4.1 by specializing the proof of that theorem to the stationary case.) As noted in the Introduction the following result is similar to the regularity condition (S) of Schäl [14]:

5.1. THEOREM. *Let the assumptions of Theorem 4.1 and also the stationarity assumptions above be satisfied. Then there exists a stationary fully optimal policy.*

Proof. Let δ^* be the fully optimal policy of Theorem 4.1. Note that δ^* need not be stationary. However, define δ^{**} by

$$\delta_k^{**}(B | x, a) = \delta_k^{**}(B | x_{k-1}) = \delta_1^*(B | x_{k-1})$$

for $B \in \mathcal{B}_k = \mathcal{B}_1$, $x_{k-1} \in X_1$. Then δ^{**} is stationary. It can readily be shown that $R(\cdot, \delta^{**}) = R(\cdot, \delta^*)$, thus implying that the stationary policy δ^{**} is fully optimal.

REFERENCES

- [1] R. Bellman, *Dynamic programming*, Princeton University Press, 1957.
- [2] D. Blackwell, *Discounted dynamic programming*, Ann. Math. Statist. 36 (1965), p. 226-235.
- [3] — *Positive dynamic programming*, Proc. Fifth Berkeley Symp. Math. Stat. Prob. 1 (1967), p. 415-418.
- [4] L. D. Brown, *Closure theorems for procedures in sequential-design processes*, p. 57-91 in: S. S. Gupta and D. S. Moore (editors), *Statistical decision theory and related topics, II*, Academic Press, 1977.
- [5] — and R. Purves, *Measurable selections of extrema*, Ann. Statist. 1 (1973), p. 902-912.
- [6] M. L. Eaton, *Complete class theorems derived from conditional complete class theorems*, ibidem 6 (1978), p. 820-827.
- [7] D. A. Freedman, *The optimal reward operator in special classes of dynamic programming problems*, Ann. Probability 2 (1974), p. 942-949.
- [8] N. Furukawa, *Markovian decision processes with compact action spaces*, Ann. Math. Statist. 43 (1972), p. 1612-1622.
- [9] K. Hinderer, *Foundations of non-stationary dynamic programming with discrete time parameter*, Springer-Verlag, 1970.
- [10] R. P. Kertz and D. C. Nachman, *Optimal plans for discrete-time non-stationary dynamic programming with general total reward function, I, II*, Tech. Repts., Georgia Institute of Technology, 1977.
- [11] L. LeCam, *An extension of Wald's theory of statistical decision functions*, Ann. Math. Statist. 26 (1955), p. 69-81.
- [12] J. Neveu, *Mathematical foundations of the calculus of probability*, Holden-Day, 1965.
- [13] M. Schäl, *On dynamic programming: Compactness of the space of policies*, Stochastic Processes Appl. 3 (1975), p. 345-364.
- [14] — *Conditions for optimality in dynamic programming and the limit of n-stage optimal policies to be optimal*, Z. Wahrscheinlichkeitstheorie und Verw. Gebiete 32 (1975), p. 179-196.
- [15] — *On dynamic programming and statistical decision theory*, Ann. Statist. 7 (1979), p. 432-445.
- [16] R. E. Strauch, *Negative dynamic programming*, Ann. Math. Statist. 37 (1966), p. 871-890.

Department of Mathematics
 Cornell University
 Ithaca, New York 14853
 White Hall, U.S.A.

Received on 1. 12. 1979