



**HAL**  
open science

## Solution methods for a class of finite-horizon vector-valued Markov decision processes

Anas Mifrani, Philippe Saint-Pierre, Nicolas Savy

► **To cite this version:**

Anas Mifrani, Philippe Saint-Pierre, Nicolas Savy. Solution methods for a class of finite-horizon vector-valued Markov decision processes. 2025. hal-04924721

**HAL Id: hal-04924721**

**<https://hal.science/hal-04924721v1>**

Preprint submitted on 31 Jan 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1                   **SOLUTION METHODS FOR A CLASS OF FINITE-HORIZON**  
2                   **VECTOR-VALUED MARKOV DECISION PROCESSES**

3                   ANAS MIFRANI\*, PHILIPPE SAINT-PIERRE, NICOLAS SAVY

ABSTRACT. This paper investigates and develops solution methods for a class of finite-horizon Markov decision processes characterized by additive or multiplicative vector rewards. Two concepts of optimality are treated: (1) optimality in the space of return vectors, whereby a policy is optimal if it delivers a maximal total reward from any initial state; and (2) optimality in the space of return functions, whereby a policy is optimal if its total reward function is maximal among all total reward functions. The paper elucidates the relation between the two concepts, proposes a procedure for utilizing this relation to determine the set of optimal policies under concept (1), and formulates a dynamic programming approach to calculating optimal policies under concept (2). The paper demonstrates that dynamic programming yields all optimal policies under concept (2). The paper’s results are illustrated with numerical experiments and a multi-objective stochastic inventory control problem.

4   *Keywords:* Multi-objective Markov decision processes; vector maximization; dynamic program-  
5   ming; multiple-criteria decision analysis.

6  
7                   1. INTRODUCTION

8       Markov decision processes offer a mathematical framework for modeling and solving sequen-  
9       tial decision making problems where outcomes are uncertain. There are three components to  
10      such a process: a stochastic dynamical system to be controlled over a period of  $N \geq 1$  epochs;  
11      real-valued rewards accrued between consecutive epochs as a result of decisions taken at epochs;  
12      and a control policy that prescribes actions such that the total expected reward (or cost) for  
13      (of) operating the system is maximized (minimized). Viewed in this way, a Markov decision  
14      process defines a single-objective, discrete-time optimal control problem.

15      However, a number of decision making problems are inherently multi-objective. In admin-  
16      istering chemotherapy, for instance, an oncologist wants to maximize the probability of cure  
17      while minimizing damage to normal cells (Coldman & Murray, 2000). A logistics manager  
18      looks for measures that simultaneously minimize the costs of warehousing and transportation

---

\*: Corresponding author.

Authors’ affiliation: Toulouse Mathematics Institute, University of Toulouse – CNRS UPS, 118 Rte de Narbonne, 31400 Toulouse, France.

Email addresses:    anas.mifrani@math.univ-toulouse.fr;    philippe.saint-pierre@math.univ-toulouse.fr;  
                          nicolas.savy@math.univ-toulouse.fr

1 over the coming year (Burns, Hall, Blumenfeld, & Daganzo, 1985). And in planning periodic  
 2 pavement rehabilitation, the local government wants to ensure the highest quality roads for  
 3 its citizens with minimal maintenance expenditures (Golabi, Kulkarni, & Way, 1982). While  
 4 the dynamics of such problems may lend themselves to Markov decision process formulation  
 5 (Puterman, 2014), it is sometimes unclear how the various objectives involved – e.g., probabil-  
 6 ity of cure versus damage to normal cells, warehousing versus transportation costs, and road  
 7 quality versus maintenance expenses – can be condensed into a single scalar-valued reward (or  
 8 cost) function (Brown & Strauch, 1965; Zadeh, 1963). Interest in overcoming this issue, and  
 9 therefore in expanding the use of Markov decision processes to multi-objective decision making,  
 10 led to the introduction of vector-valued Markov decision processes.

11 In a vector-valued Markov decision process, rewards take values in  $\mathbb{R}^m$ ,  $m \geq 1$ , with each  
 12 reward component representing an optimization objective. The standard formulation is as  
 13 follows. Let  $S$  denote the set of states the system can occupy throughout its lifetime. For any  
 14 state  $s$ , let  $A_s$  be the set of actions available in  $s$ ;  $R_t(s, a) = (r_t(s, a)_1, \dots, r_t(s, a)_m)$  the reward  
 15 for choosing  $a \in A_s$  in  $s$  at time  $t = 1, \dots, N - 1$ , and  $R_N(s)$  the reward for occupying state  
 16  $s$  at the terminal epoch;  $p_t(j|s, a)$  the transition probability from  $s$  to  $j \in S$  if  $a \in A_s$  was  
 17 chosen at time  $t$ ; and  $u_t^\pi(s) = \mathbb{E}_\pi^s[\sum_{i=t}^{N-1} R_i(X_i, d_i(X_i)) + R_N(X_N)] \in \mathbb{R}^m$  the expected total  
 18 reward for using a Markovian deterministic policy  $\pi$  from  $t$  onward given  $X_t = s$ , where  $X_i$   
 19 represents the (random) state at time  $i$  and  $d_i(X_i)$  the action prescribed by  $\pi$  for  $X_i$  at time  
 20  $i$ . The last expectation is called a policy return. In particular, the vector  $u_1^\pi(s)$  represents the  
 21 return of a policy  $\pi$  over the entire decision making horizon. When  $m = 1$ , this model reduces  
 22 to a Markov decision process.

23 Roughly speaking, the reward structure just described induces a partial order on the set  
 24 of policies whereby a policy  $\pi$  may be superior to a policy  $\pi'$  in some respects but inferior  
 25 to  $\pi'$  in others. For example, taking  $m = 2$  and a common initial state  $s$ , we may have  
 26  $u_1^\pi(s)_1 \geq u_1^{\pi'}(s)_1$  yet  $u_1^\pi(s)_2 < u_1^{\pi'}(s)_2$ , so that, componentwise, we neither have  $u_1^\pi(s) \geq u_1^{\pi'}(s)$   
 27 nor  $u_1^{\pi'}(s) \geq u_1^\pi(s)$ . In a standard Markov decision process ( $m = 1$ ), this situation obviously  
 28 never arises.

29 Though Brown and Strauch (1965) were the first to consider a Markov decision process  
 30 with partially ordered rewards, namely rewards in multiplicative lattices, the chief theoret-  
 31 ical developments concerning vector-valued Markov decision processes as presented here oc-  
 32 curred in papers published between the 1970s and 1980s (Furukawa, 1980; Henig, 1983; White,  
 33 1982). To the best of our knowledge, D. J. White’s seminal paper (White, 1982) contains the  
 34 first attempt at formulating an exact dynamic programming approach to solving a class of  
 35 vector-valued Markov decision processes. This approach has been cited by a recent survey of  
 36 multi-objective reinforcement learning (Hayes et al., 2022), and numerous authors have used  
 37 it either to compare it experimentally with their own approaches (Roijers, Röpke, Nowe, &

1 Radulescu, 2021; Wiering & De Jong, 2007) or as a point of departure for the development of  
 2 new algorithms (Chen, Trevizan, & Thiébaux, 2023; Mandow, Pérez-de-la Cruz, & Pozas, 2022;  
 3 Ruiz-Montiel, Mandow, & Pérez-de-la Cruz, 2017; Van Moffaert & Nowé, 2014). It is based  
 4 on the following vector analogue of the Bellman equations of a finite-horizon Markov decision  
 5 process (Puterman, 2014, Chapter 4):

$$U_t(s) = e \left( \bigcup_{a \in A_s} \left( \{R_t(s, a)\} \oplus \sum_{j \in S} p_t(j|s, a) U_{t+1}(j) \right) \right); \quad t < N \quad (1)$$

$$U_t(s) = \{R_N(s)\}; \quad t = N \quad (2)$$

6  
 7 for all  $s \in S$  and  $t = 1, \dots, N$ , where  $e(X)$  denotes the Pareto efficient subset of a set  $X \subseteq \mathbb{R}^m$   
 8 (see Section 2 for a formal definition),  $A \oplus B = \{a + b : \forall a \in A, \forall b \in B\}$  for any two nonempty  
 9 sets  $A$  and  $B$ , and where the unknowns are the  $U_t(s)$ 's,  $s \in S$ ,  $t = 1, \dots, N$ .

10 White claims that the solutions of Equations (1) and (2) are the Pareto efficient sets of policy  
 11 returns for all epochs and initial states, i.e.  $U_t(s) = e(\bigcup_{\pi} \{u_t^{\pi}(s)\})$  for all  $t \leq N$  and  $s \in S$   
 12 (White, 1982, Theorem 2). In fact this claim is generally false (Mifrani, 2023), notwithstanding  
 13 its coincidence, for  $m = 1$ , with the correct observation that the solutions of the Bellman  
 14 equations are the  $\max_{\pi} u_t^{\pi}(s)$ 's (Puterman, 2014, Proposition 4.3.3.).

15 We might note, in passing, that such issues do not arise in infinite horizon models ( $N = \infty$ ).  
 16 Furukawa (1980) has proved that the fixed-point characterization of a Markov decision process's  
 17 optimal infinite horizon value (Puterman, 2014, Theorem 6.2.6.) extends *mutatis mutandis* to  
 18 vector-valued processes. In short, the infinite horizon counterparts of White's equations are  
 19 valid. Here we shall confine our analysis to finite horizon models.

20 In this paper, we take the position of a decision maker who has to select a *V-optimal* (V  
 21 for vector-based) Markovian deterministic policy, that is, a policy which generates an efficient  
 22 return from any initial state. We develop an approach to computing such a policy that does  
 23 not involve dynamic programming on the space of return vectors. The key role in this approach  
 24 is played by an auxiliary optimality criterion that we call *F-optimality* (F for function-based).  
 25 The difference between the two criteria lies in that, to compare a pair of policies  $\pi$  and  $\pi'$ ,  
 26 F-optimality focuses on the policies' return *functions*,  $u_t^{\pi}$  and  $u_t^{\pi'}$ ,  $t = 1, \dots, N$ , rather than on  
 27 the return *vectors*  $u_t^{\pi}(s)$ ,  $u_t^{\pi'}(s)$  achieved in individual states  $s$ . The subtlety of this distinction  
 28 will be illustrated in Example 1 of Section 3.

29 We establish the following: (1) V-optimality is subsumed under F-optimality; (2) F-optimality  
 30 is susceptible to dynamic programming; (3) the solutions to the dynamic programming equa-  
 31 tions can be leveraged to construct F-optimal policies; and (4) provided there is a finite number  
 32 of F-optimal policies, a computationally useful characterization of V-optimal policies within the  
 33 set of F-optimal policies can be implemented to find all policies of the former kind. Thus, in  
 34 particular, we shall see that all V-optimal policies can be calculated without evaluating the

1  $e(\bigcup_{\pi}\{u_1^{\pi}(s)\})$ 's, a potentially intractable task in the absence of a valid recurrence relation  
 2 between  $e(\bigcup_{\pi}\{u_t^{\pi}(s)\})$  and  $e(\bigcup_{\pi}\{u_{t+1}^{\pi}(j)\})$  for all  $s, j \in S$  and  $t = 1, \dots, N - 1$ .

3 The hypotheses and notation underpinning this paper are presented in greater detail in Sec-  
 4 tion 2. In Section 3, we shall substantiate, and discuss the implications of, points (1)-(4) as  
 5 outlined above. In particular, we shall devise algorithms for computing policies according to  
 6 each criterion. A numerical analysis of the algorithms is undertaken in Section 3. Section 4  
 7 reports implementation results for a multi-objective stochastic inventory management prob-  
 8 lem. In Section 5, we consider the ramifications of our results for models with multiplicative  
 9 – rather than additive – rewards, make some general comments on the algorithms, and close  
 10 with a discussion of potential applications of these results.

11

12

## 2. MODEL ASSUMPTIONS AND NOTATION

13 At each epoch  $t \leq N$ , the system occupies a state  $s_t$ . The set of all states,  $S$ , is finite. The  
 14 decision maker has at their disposal a set of actions,  $A$ , which they must choose from at each  
 15 epoch. If only certain actions are allowed in a state, let  $A_s$  be the set of permissible actions in  
 16  $s \in S$ , from which it follows  $A = \bigcup_{s \in S} A_s$ . Suppose  $A_s$  is a compact subset of  $\mathbb{R}$  for all  $s \in S$ .  
 17 Assuming  $a \in A$  was selected at time  $t < N$ , the probability that the system will occupy state  
 18  $j \in S$  at  $t + 1$  depends only on the present state  $s \in S$ , and is denoted by  $p_t(j|s, a)$ . For  
 19 choosing action  $a \in A$  in state  $s$  at time  $t < N$ , the decision maker receives a vector reward  
 20  $R_t(s, a) \in \mathbb{R}^m$ ,  $m \geq 2$ . Suppose that transition probabilities and rewards are continuous on  
 21  $A_s$  for all  $s \in S$ . A (Markovian, deterministic) decision rule  $d_t$  dictates the action to be taken  
 22 in each state at epoch  $t < N$ , and is viewed therefore as a mapping from  $S$  to  $A$ . The set  
 23 of all decision rules,  $D$ , is considered to be compact. For any  $t < N$  and any  $d_t \in D$ , let  
 24  $P^{d_t} = (p_t(j|s, d_t(s)))_{s, j \in S}$  be the transition probability matrix induced by  $d_t$ . No decision is  
 25 taken at epoch  $N$ , but a state-dependent reward  $R_N(s)$  is generated. A policy specifies the  
 26 decision rule that should be used at each epoch, and shall be identified with its corresponding  
 27 sequence of decision rules  $(d_1, \dots, d_{N-1})$ . Let  $\Pi = D^{N-1}$  be the set of all policies. For any  
 28  $\pi \in \Pi$  and  $t < N$ ,  $\bar{\pi}(t) = (d_t, \dots, d_{N-1})$  shall denote the portion of decision rules used by  $\pi$   
 29 from  $t$  onward.

30 For any policy  $\pi = (d_1, \dots, d_{N-1})$  and any  $t < N$ , we have the recurrence relation

$$u_t^{\pi}(s) = R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s)) u_{t+1}^{\pi}(j) \quad (3)$$

31 where we let

$$u_N^{\pi}(s) = R_N(s). \quad (4)$$

1 Expanding the sum in (3) over all future epochs and states yields the expression

$$u_t^\pi(s) = R_t(s, d_t(s)) + \sum_{i=t}^{N-2} \sum_{j \in S} \left( \prod_{k=t}^i P^{d_k} \right)_{s,j} R_{i+1}(j, d_{i+1}(j)) + \sum_{j \in S} \left( \prod_{k=t}^{N-1} P^{d_k} \right)_{s,j} R_N(j). \quad (5)$$

2 The terms “policy return” and “return”, where time and state are omitted for brevity, shall  
 3 refer to any vector  $u \in \mathbb{R}^m$  for which there is a policy  $\pi$ , a time  $t = 1, \dots, N$  and an  $s \in S$  with  
 4  $u = u_t^\pi(s)$ . Where a distinction must be drawn between the function  $u_t^\pi$ , and the values  $u_t^\pi(s)$   
 5 it takes at particular states  $s \in S$ , the phrases “(policy) return function” and “(policy) return  
 6 vector” shall be used instead, with time and state also omitted for brevity.

7 For any partially ordered set  $(X, \geq)$ , let  $e(X)$  be the efficient (or admissible, or noninferior,  
 8 or Pareto optimal) subset of  $X$ , to wit:

$$e(X) = \{x \in X : \forall y \in X, y \geq x \implies y = x\}. \quad (6)$$

9 Let  $F(S, \mathbb{R}^m)$  denote the set of all  $\mathbb{R}^m$ -valued functions on  $S$ . In Section 3 we shall be concerned  
 10 with efficiency in subsets of  $(\mathbb{R}^m, \geq)$  and  $(F(S, \mathbb{R}^m), \succeq)$ , where:

$$\forall x, y \in \mathbb{R}^m, x \geq y \iff \forall i = 1, \dots, m, x_i \geq y_i, \quad (7)$$

$$\forall u, v \in F(S, \mathbb{R}^m), u \succeq v \iff \forall s \in S, u(s) \geq v(s). \quad (8)$$

12 The partial orders thus defined provide a means for comparing, respectively, return vectors  
 13 and return functions. A strict partial order  $>$  can also be defined on  $(X, \geq)$  as  $\forall x, y \in X, x >$   
 14  $y \iff x \geq y \wedge x \neq y$ . When  $X \subseteq \mathbb{R}^m$  and  $\mathbb{R}^m$  is equipped with (7), the elements of  $e(X)$   
 15 are sometimes referred to as “vector maxima” (Geoffrion, 1968), though for consistency with  
 16 previous work on vector-valued Markov decision processes the generic adjective “efficient” shall  
 17 be used instead. When  $X$  has a maximum, such as is the case with  $\bigcup_{\pi \in \Pi} \{u_t^\pi(s)\}$  for  $m = 1$   
 18 (Puterman, 2014, Proposition 4.4.3), we have  $e(X) = \{\max(X)\}$ .

19

20

### 3. THEORETICAL RESULTS

21 As stated in the Introduction, we shall study two related concepts of optimality as regards  
 22 policies. In the first concept, a policy is optimal if, whatever the state in which it was first  
 23 implemented, it delivers a maximal return over the  $N$  epochs:

24 **Definition 1** (V-optimality). *A policy  $\pi^V$  is V-optimal if and only if  $u_1^{\pi^V}(s) \in e(\bigcup_{\pi} \{u_1^\pi(s)\})$   
 25 for all states  $s \in S$ .*

26 Here “V” stands for “vector”, and  $\bigcup_{\pi} \{u_1^\pi(s)\}$  is endowed with  $\geq$  as defined in (7). For  $m = 1$ ,  
 27  $\max_{\pi} u_1^\pi(s)$  exists for all  $s \in S$  (Puterman, 2014, Proposition 4.3.3.), and Definition 1 reads  
 28 “ $\pi^V$  is V-optimal if and only if  $u_1^{\pi^V}(s) = \max_{\pi} u_1^\pi(s)$  for all  $s \in S$ ”, which is the standard  
 29 optimality criterion across a wide range of Markov decision process applications (Borrero &  
 30 Akhavan-Tabatabaei, 2013; Goedhart, Haijema, Akkerman, & de Leeuw, 2023; Mason, Denton,

1 [Shah, & Smith, 2014](#); [Puterman, 2014](#); [Ramirez-Nafarrate, Hafizoglu, Gel, & Fowler, 2014](#);  
2 [Schlosser & Gönsch, 2023](#)).

3 One of the aims of this section is to supply a procedure for determining all V-optimal  
4 policies under the hypotheses of Section 2. This will be achieved by leveraging the connection  
5 between V-optimality and a neighboring optimality concept of which return functions, rather  
6 than return vectors, are the core ingredient.

7 **Definition 2** (F-optimality). *A policy  $\pi^F$  is F-optimal if and only if  $u_1^{\pi^F} \in e(\bigcup_{\pi \in \Pi} \{u_1^\pi\})$ ,*

8 where “F” stands for “function”, and where it is implicit that  $\bigcup_{\pi \in \Pi} \{u_1^\pi\}$  is ordered by  $\succeq$  as  
9 defined in (8).

10 In brief, a policy  $\pi^*$  is V-optimal if for each state  $s$  there exists no other policy  $\pi_s \neq \pi^*$  with  
11  $u_1^{\pi_s}(s) \geq u_1^{\pi^*}(s)$ , and is F-optimal if there is no other  $\pi$  such that  $u_1^\pi \succeq u_1^{\pi^*}$ .

12 **Example 1.** *The distinction is illustrated by the following situation. Suppose these policies*  
13 *were available in a two-state model with  $m = 2$ : a policy  $\pi_1$  yielding  $u_1^{\pi_1}(s_1) = (3, 1)$  and*  
14  *$u_1^{\pi_1}(s_2) = (5, -2)$ , and a policy  $\pi_2$  yielding  $u_1^{\pi_2}(s_1) = (2, 1)$  and  $u_1^{\pi_2}(s_2) = (\frac{1}{2}, 0)$ . Then  $u_1^{\pi_1}$  and*  
15  *$u_1^{\pi_2}$  are incomparable with respect to  $\succeq$ ;  $u_1^{\pi_1}(s_2)$  and  $u_1^{\pi_2}(s_2)$  are incomparable with respect to*  
16  *$\geq$ ; and  $u_1^{\pi_1}(s_1) > u_1^{\pi_2}(s_1)$ . By definition,  $\pi_2$  is not V-optimal, for  $u_1^{\pi_1}(s_1) > u_1^{\pi_2}(s_1)$  implies the*  
17 *existence of a state  $s$  ( $s_1$  here) for which there is a policy  $\pi_s \neq \pi_2$  ( $\pi_1$  here) with  $u_1^{\pi_s}(s) \geq u_1^{\pi_2}(s)$ .*  
18 *It may still be F-optimal, however, as we have  $u_1^{\pi_1}(s_2) \not\geq u_1^{\pi_2}(s_2)$  and therefore  $u_1^{\pi_1} \not\geq u_1^{\pi_2}$ . If*  
19 *some third policy  $\pi_3$  satisfied  $u_1^{\pi_3}(s_1) \geq u_1^{\pi_2}(s_1)$  and  $u_1^{\pi_3}(s_2) \geq u_1^{\pi_2}(s_2)$ , then  $\pi_2$  would not be*  
20 *F-optimal.*

21 In the succeeding development, we will find it convenient to focus on the latter concept for  
22 four key reasons, all of which will be demonstrated in due course: (1) we are able to guarantee  
23 the existence of F-optimal policies; (2) the problem of finding F-optimal policies is susceptible  
24 to dynamic programming; (3) F-optimal policies satisfy the Principle of Optimality ([Bellman,](#)  
25 [1954](#)); and (4) a V-optimal policy *must* be F-optimal, that is, given  $\pi \in \Pi$ , efficiency of  $u_1^\pi$   
26 in  $F(S, \mathbb{R}^m)$  is a necessary condition for efficiency of  $u_1^\pi(s)$  in  $\mathbb{R}^m$  for all  $s \in S$ . These obser-  
27 vations have important practical implications. First, the fact that the Principle of Optimality  
28 holds means that all F-optimal – and therefore all V-optimal – policies will be found through  
29 dynamic programming. Second, if we can determine which policies are not F-optimal, we will  
30 immediately recognize those that are not V-optimal. Third, if  $e(\bigcup_{\pi \in \Pi} \{u_1^\pi\})$  is finite and can be  
31 computed in a finite number of steps, it will be possible to obtain the set of V-optimal policies,  
32 thereby solving both optimization problems simultaneously. This last point will be illustrated  
33 in the next section.

34 **Proposition 1.** *Let  $\pi^* \in \Pi$ . If  $\pi^*$  is V-optimal, then it is F-optimal.*

1 *Proof.* Suppose that  $\pi^* \in \Pi$  is V-optimal. Let  $\pi' \in \Pi$  be a policy such that  $u_1^{\pi'} \succeq u_1^{\pi^*}$ . Then for  
 2 any  $s \in S$ ,  $u_1^{\pi'}(s) \geq u_1^{\pi^*}(s)$ . Therefore, since  $\pi^*$  is V-optimal, it follows that  $u_1^{\pi'}(s) = u_1^{\pi^*}(s)$   
 3 for any  $s \in S$ . Thus,  $u_1^{\pi'} = u_1^{\pi^*}$ . This shows that  $u_1^{\pi^*} \in e(\bigcup_{\pi} \{u_1^{\pi}\})$ , and hence that  $\pi^*$  is  
 4 F-optimal.  $\square$

5 Lemma 1 formalizes a useful intuition about return functions that will be invoked repeatedly  
 6 throughout this section.

7 **Lemma 1.** *Let  $\pi = (d_1, \dots, d_{N-1})$ ,  $\pi' = (d'_1, \dots, d'_{N-1}) \in \Pi$ ,  $d_t \in D$  and  $t = 1, \dots, N-2$ . Suppose  
 8  $u_{t+1}^{\pi} \succeq u_{t+1}^{\pi'}$ . Then for any two policies  $\pi_1$  and  $\pi_2$  such that  $\bar{\pi}_1(t) = (d_t, d_{t+1}, \dots, d_{N-1})$  and  
 9  $\bar{\pi}_2(t) = (d_t, d'_{t+1}, \dots, d'_{N-1})$ ,  $u_t^{\pi_1} \succeq u_t^{\pi_2}$ .*

10 *Proof.* Suppose  $u_{t+1}^{\pi} \succeq u_{t+1}^{\pi'}$ , and let  $s \in S$ . Let  $\pi_1, \pi_2 \in \Pi$  be policies such that  $\bar{\pi}_1(t) =$   
 11  $(d_t, d_{t+1}, \dots, d_{N-1})$  and  $\bar{\pi}_2(t) = (d_t, d'_{t+1}, \dots, d'_{N-1})$ . For all  $j \in S$ ,  $u_{t+1}^{\pi_1}(j) = u_{t+1}^{\pi}(j) \geq$   
 12  $u_{t+1}^{\pi'}(j) = u_{t+1}^{\pi_2}(j)$ , hence  $\sum_{j \in S} p(j|s, d_t(s))u_{t+1}^{\pi_1}(j) \geq \sum_{j \in S} p(j|s, d_t(s))u_{t+1}^{\pi_2}(j)$  due to the  
 13 nonnegativity of probabilities. Thus,

$$R_t(s, d_t(s)) + \sum_{j \in S} p(j|s, d_t(s))u_{t+1}^{\pi_1}(j) \geq R_t(s, d_t(s)) + \sum_{j \in S} p(j|s, d_t(s))u_{t+1}^{\pi_2}(j).$$

14 This establishes  $u_t^{\pi_1}(s) \geq u_t^{\pi_2}(s)$  for each  $s \in S$ . Ergo,  $u_t^{\pi_1} \succeq u_t^{\pi_2}$ .  $\square$

15 **Example 2.** *Consider a vector-valued Markov decision process with  $S = \{1, 2\}$ ,  $A_1 = \{a, b\}$ ,  
 16 and  $A_2 = \{a\}$ . Suppose that at a decision epoch  $t$  we had  $p_t(1|1, a) = .75$ ;  $p_t(2|1, a) = .25$ ;  
 17  $p_t(1|1, b) = p_t(2|1, b) = .5$ ;  $p_t(1|2, a) = 1$ ;  $p_t(2|2, a) = 0$ ;  $R_t(1, a) = (1, 0)$ ;  $R_t(2, a) = (0, 0)$ ;  
 18 and  $R_t(1, b) = (0, 1)$ .*

19 *For the purposes of this example, let us assume that there exist policies  $\pi$  and  $\pi'$  with  
 20 returns  $u_{t+1}^{\pi}(1) = (0, 0)$ ,  $u_{t+1}^{\pi}(2) = (-2, 2)$ ,  $u_{t+1}^{\pi'}(1) = (-0.5, 0)$ ,  $u_{t+1}^{\pi'}(2) = (-6, 1)$ . Clearly,  
 21  $u_{t+1}^{\pi} \succeq u_{t+1}^{\pi'}$ .*

22 *Now let  $d_t \in D$  denote the decision rule that chooses  $b$  in state 1, i.e  $d_t(1) = b$  and  $d_t(2) = a$ .  
 23 Choose  $\pi_1$  to be any policy that selects  $d_t$  at time  $t$  then pursues  $\pi$  from time  $t+1$  onward.  
 24 Similarly, let  $\pi_2$  select  $d_t$  at time  $t$  then pursue  $\pi'$  at all future epochs. In our notation,  
 25 this translates to  $\bar{\pi}_1(t) = (d_t, \pi)$  and  $\bar{\pi}_2(t) = (d_t, \pi')$ . Through simple calculations, we will  
 26 demonstrate the assertion in Lemma 1 that  $u_t^{\pi_1} \succeq u_t^{\pi_2}$ . From Equation (3) we have that*

$$\begin{aligned} u_t^{\pi_1}(1) &= R_t(1, d_t(1)) + \sum_{j \in S} p_t(j|1, d_t(1))u_{t+1}^{\pi}(j) \\ &= (0, 1) + 0.5 \cdot (0, 0) + 0.5 \cdot (-2, 2) \\ &= (-1, 2), \end{aligned}$$



1 and

$$\begin{aligned}
 u_t^{\pi_2}(1) &= R_t(1, d_t(1)) + \sum_{j \in S} p_t(j|1, d_t(1)) u_{t+1}^{\pi_2'}(j) \\
 &= (0, 1) + 0.5 \cdot (-0.5, 0) + 0.5 \cdot (-6, 1) \\
 &= (-3.25, 1.5).
 \end{aligned}$$

2 Thus,  $u_t^{\pi_1}(1) \geq u_t^{\pi_2}(1)$ . The reader can easily replicate this method of calculation to verify that  
 3  $u_t^{\pi_1}(2) = (0, 0)$  and  $u_t^{\pi_2}(2) = (-0.5, 0)$ . This means that  $u_t^{\pi_1}(2) \geq u_t^{\pi_2}(2)$ , hence  $u_t^{\pi_1} \succeq u_t^{\pi_2}$ .

4 Fundamental to the proof of Lemma 1 is the fact that for any  $\pi = (d_1, \dots, d_{N-1}) \in \Pi$ ,  $s \in S$   
 5 and  $t = 1, \dots, N-1$ ,  $u_t^\pi(s) = R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s)) u_{t+1}^\pi(j)$ . That is, policy returns  
 6 are separable and additive. The scope of the lemma, however, covers a broader category of  
 7 separable returns. Following [Morin \(1982\)](#), we can make this generalization: any vector-valued  
 8 Markov decision process such that  $u_t^\pi(s) = R_t(s, d_t(s)) \circ \sum_{j \in S} p_t(j|s, d_t(s)) u_{t+1}^\pi(j)$ , where  
 9  $\circ$  is an isotonic symmetric binary operator, i.e a symmetric binary operator that preserves  
 10 inequalities (with respect to  $\geq$ ), satisfies the lemma. To prove this generalization, we may  
 11 proceed in exactly the same fashion as above, concluding from the isotonicity of  $\circ$  that

$$R_t(s, d(s)) \circ \sum_{j \in S} p(j|s, d(s)) u_{t+1}^{\pi_1}(j) \geq R_t(s, d(s)) \circ \sum_{j \in S} p(j|s, d(s)) u_{t+1}^{\pi_2}(j)$$

12 for all  $s \in S$ , and therefore that  $u_t^{\pi_1} \succeq u_t^{\pi_2}$ . Incidentally, Morin points out that a strictly  
 13 isotonic associative  $\circ$ , of which addition in  $\mathbb{R}^m$  and componentwise multiplication in  $(0, \infty)^m$   
 14 would be examples, ensures the validity of the Bellman equations in Markov decision processes.  
 15 However, [Mifrani \(2023\)](#) has recently shown that this is not true for all vector-valued Markov  
 16 decision processes with regard to the vector extension of those equations.

17 It will later prove desirable, especially for the purpose of justifying optimality equations, to  
 18 have a property that enables us to assert that each inefficient point in  $\bigcup_{\pi \in \Pi} \{u_t^\pi\}$  is dominated  
 19 by an efficient one. Notice that this is not entailed by the definition of efficiency, because in  
 20 general, all we can say about an inefficient point is that it is dominated by another point, which  
 21 may or may not be efficient. [Berge \(1985\)](#) calls “absorbent” a partially ordered set  $S \subseteq (X, \geq)$   
 22 such that for every  $x \in X$ , there exists  $s \in S$  satisfying  $s \geq x$ . We wish then an absorbent  
 23  $e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$  for all  $t = 1, \dots, N$ , so that in addition to the aforementioned property, we may  
 24 conclude that  $e(\bigcup_{\pi \in \Pi} \{u_1^\pi\}) \neq \emptyset$ , and therefore that F-optimal policies exist. The following  
 25 lemma from [Henig \(1985\)](#) implies that a nonempty partially ordered set is absorbent if it meets  
 26 the conditions of Zorn’s lemma ([Zorn, 1935](#)).

27 **Lemma 2.** ([Henig, 1985](#)) Let  $(U, \geq)$  be a nonempty partially ordered set, and  $K$  a nonempty  
 28 subset of  $U$ . Suppose that for every  $u \in U$  and every  $v \in K$ ,  $u \geq v$  implies  $u \in K$ . Suppose  
 29 further that every totally ordered subset (chain) of  $U$  has an upper bound in  $U$ . Then  $e(U) \cap K \neq$   
 30  $\emptyset$ .

1 For fixed  $u \in \bigcup_{\pi \in \Pi} \{u_t^\pi\}$  and  $t = 1, \dots, N$ , let  $K(u)$  denote the set  $\{v \in \bigcup_{\pi} \{u_t^\pi\} : v \succeq u\}$ . We  
 2 claim that there is an efficient  $v \in e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$  such that  $v \succeq u$ . Our proof rests on Lemma  
 3 2. First of all,  $K(u)$  is nonempty, as  $u \succeq u$ . Moreover, as will be shown below,

- 4 (1) for every  $v$  in  $K(u)$  and every  $v' \in \bigcup_{\pi} \{u_t^\pi\}$ ,  $v' \succeq v$  implies  $v' \in K(u)$ , and;  
 5 (2) every chain of  $\bigcup_{\pi \in \Pi} \{u_t^\pi\}$  is bounded above in  $\bigcup_{\pi \in \Pi} \{u_t^\pi\}$ .

6 The proof of point (2) involves studying the convergence of certain sequences in  $D$  and in  
 7  $\Pi = D^{N-1}$ . For convergence to be meaningful on either set, a topology must be introduced. The  
 8 most widely assumed topology in the analysis of Markov decision processes is that of uniform  
 9 (or sup-norm) convergence. But because uniform convergence implies pointwise convergence,  
 10 and because our proof does not use properties of the former which are not true of the latter, it  
 11 suffices to endow  $D$  with the topology of pointwise convergence and, by extension,  $\Pi$  with the  
 12 associated product topology.

13 **Theorem 1.** *Equip  $D$  with the topology of pointwise convergence and  $\Pi = D^{N-1}$  with the*  
 14 *product topology. Then points (1) and (2) as enunciated above are true.*

15 *Proof.* We divide the proof into two parts.

16 (1) Let  $v \in K(u)$  and  $v' \in \bigcup_{\pi} \{u_t^\pi\}$ . If  $v' \succeq v$ , then, since  $v \succeq u$  and  $\succeq$  is transitive, we have  
 17 that  $v' \succeq u$ , hence  $v' \in K(u)$ .

18 (2) Notice first that  $\Pi$ , being the product of compact sets, is compact. For all  $\pi \in \Pi$ , let  
 19  $f_t(\pi) = u_t^\pi$ . According to (Birrkhoff, 1940, Theorem 16), it suffices to show that  $f_t$ , viewed as  
 20 a mapping from  $\Pi$  to  $F(S, \mathbb{R}^m)$ , satisfies the following property: whenever  $e \in F(S, \mathbb{R}^m)$  and  
 21 for every sequence  $(\pi_n)_n$  with values in  $\Pi$ ,  $\pi_n \rightarrow \pi^\circ$  and  $f_t(\pi_n) \geq e$  for all  $n$  imply  $f_t(\pi^\circ) \geq e$ .

22 Accordingly, let  $e \in F(S, \mathbb{R}^m)$  and  $(\pi_n)_n$  a sequence of policies converging to a  $\pi^\circ \in \Pi$ , with  
 23  $f_t(\pi_n) \succeq e$  for all  $n$ . Let  $s \in S$ . From (5), we have that

$$f_t(\pi_n)(s) = R_t(s, d_t^{\pi_n}(s)) + \sum_{i=t}^{N-2} \sum_{j \in S} \left( \prod_{k=t}^i P_k^{d_k^{\pi_n}} \right)_{s,j} R_{i+1}(j, d_{i+1}^{\pi_n}(j))$$

$$+ \sum_{j \in S} \left( \prod_{k=t}^{N-1} P_k^{d_k^{\pi_n}} \right)_{s,j} R_N(j) \geq e(s)$$

24 for all  $n$ . Now, in view of  $D$ 's topology, we have that for all  $i = 1, \dots, N$ ,  $d_i^{\pi_n}(s) \rightarrow d_i^{\pi^\circ}(s)$  in  $A$ .  
 25 This, together with the continuity of the transition probabilities and of each reward component  
 26 on  $A$ , yields  $f_t(\pi_n)(s)_p \rightarrow f_t(\pi^\circ)(s)_p$  in  $\mathbb{R}$  and hence  $f_t(\pi^\circ)(s)_p \geq e(s)_p$  for all  $p = 1, \dots, m$ .  
 27 Thus, by definition,  $f_t(\pi^\circ)(s) \geq e(s)$ . Since  $s$  was chosen arbitrarily, it follows that  $f_t(\pi^\circ) \succeq e$ ,  
 28 again by definition of  $\succeq$ . By virtue of this and the compactness of  $\Pi$ , it follows from Theorem  
 29 1 that every chain in  $f_t(P) = \bigcup_{\pi} \{u_t^\pi\}$  has an upper bound in  $\bigcup_{\pi} \{u_t^\pi\}$ .  $\square$

1 Theorem 1 relies on the fact that  $\Pi$  is compact, which in turn relies on the fact that  $D$  is  
 2 compact. Considering that many Markov decision process applications use a finite  $A$  (Borrero  
 3 & Akhavan-Tabatabaei, 2013; Goedhart et al., 2023; Mason et al., 2014; Ramirez-Nafarrate et  
 4 al., 2014; Schlosser & Gönsch, 2023; Wang, Demeulemeester, Vansteenkiste, & Rademakers,  
 5 2024; White, 1993), and therefore a compact  $D$ , this is not as restrictive an assumption as it  
 6 may seem at first glance.

7 **Corollary 1.** *For all  $t = 1, \dots, N$ , if  $u \in \bigcup_{\pi \in \Pi} \{u_t^\pi\}$ , there is an efficient return function  
 8  $v \in e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$  such that  $v \succeq u$ .*

9 **Example 3.** *We illustrate part (2) of Theorem 1 in tandem with Corollary 1. In a certain  
 10 model with state space  $S = \{1, 2\}$  and  $m = 2$ , the return functions generated by all policies  
 11 from time  $t = 1$  onward were found to be given by*

$$\bigcup_{\pi \in \Pi} \{u_1^\pi\} = \{u^1, u^2, u^3, u^4, u^5, u^6, u^7, u^8\} \subset F(S, \mathbb{R}^2),$$

12 where, for example,  $u^1(1) = (1.64, 12.6)$ ;  $u^1(2) = (2.44, 8.56)$ ;  $u^4(1) = (3.44, -2.44)$ ;  $u^4(2) =$   
 13  $(1.62, -7.62)$ ;  $u^8(1) = (-4.62, -6.37)$ ;  $u^8(2) = (-3.75, -10.25)$ . The exact values are immate-  
 14 rial to the purposes of this example; what is important here are the relations among the points  
 15 in the above set. We can see that  $u^1 \succeq u^8$ ,  $u^4 \succeq u^8$ , and that no comparison is possible be-  
 16 tween  $u^1$  and  $u^4$ . The full network of relations is summarized in the diagram of Figure 1. For  
 17 example, the diagram indicates that  $u^1 \succeq u^3$ , but also that  $u^1 \succeq u^7$ , as  $u^3 \succeq u^7$  and  $\succeq$  is tran-  
 18 sitive. Such drawings are known in set theory as Hasse diagrams, and are particularly useful  
 19 for determining chains and antichains (subsets of which no distinct points are comparable) in  
 20 a partially ordered set.

21 We shall verify the assertion in Theorem 1(2) that every totally ordered subset of  $\bigcup_{\pi \in \Pi} \{u_1^\pi\}$   
 22 is upper bounded in  $\bigcup_{\pi \in \Pi} \{u_1^\pi\}$  relative to  $\succeq$ . According to Figure 1, the totally ordered subsets  
 23 in this example comprise:

- 24 (1) eight singletons  $\{u^i\}$ ,  $i = 1, \dots, 8$ , which are bounded above by virtue of  $\succeq$  being reflexive;  
 25 (2) fifteen two-point sets including  $\{u^1, u^5\}$ ,  $\{u^4, u^8\}$ ,  $\{u^1, u^8\}$ , and  $\{u^2, u^8\}$ , all of which  
 26 are bounded above by the element from which the arrow (or arrows) originates (origi-  
 27 nate);  
 28 (3) and six three-point sets  $\{u^1, u^3, u^8\}$ ,  $\{u^1, u^3, u^7\}$ ,  $\{u^2, u^4, u^8\}$ ,  $\{u^2, u^3, u^8\}$ ,  $\{u^2, u^4, u^7\}$ ,  
 29 and  $\{u^2, u^3, u^7\}$ , all of which are bounded above either by  $u^1$  or by  $u^2$ .

30 On the same diagram we can observe that  $e(\bigcup_{\pi \in \Pi} \{u_1^\pi\}) = \{u^1, u^2\}$ , since  $u^1$  and  $u^2$  are the  
 31 only points towards which no arrows are directed. If the claim in Corollary 1 is correct, then  
 32 each  $u^i$ ,  $i = 1, \dots, 8$ , should satisfy  $u^1 \succeq u^i$ , or  $u^2 \succeq u^i$ , or both. A simple inspection of Figure  
 33 1 reveals that this is indeed the case.

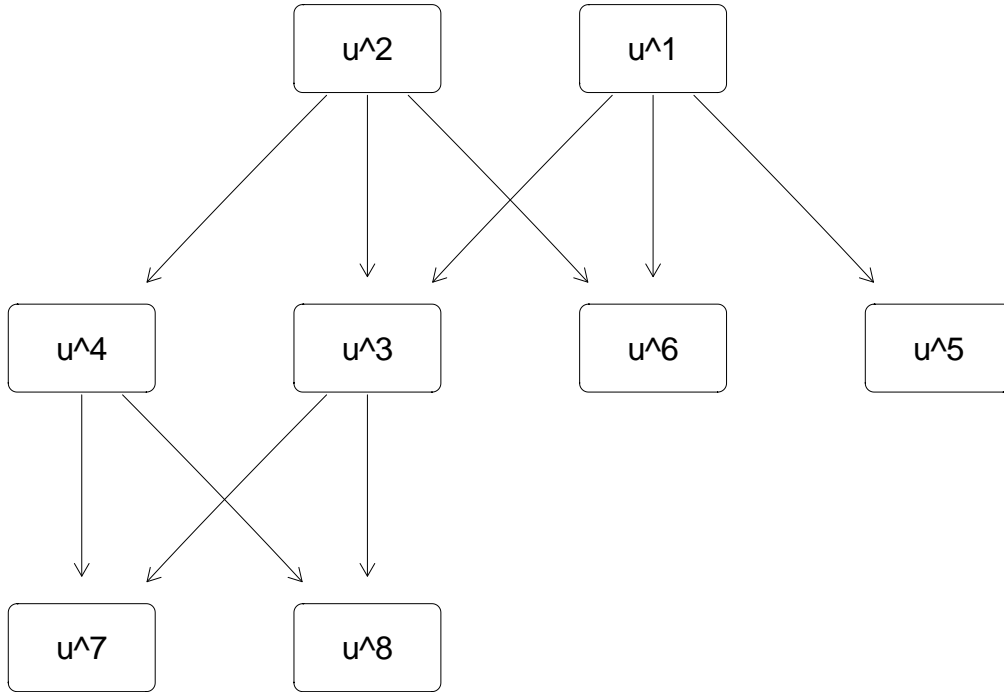


FIGURE 1. A Hasse diagram of the set  $\bigcup_{\pi \in \Pi} \{u_1^\pi\}$  in Example 3, ordered by  $\succeq$ . An outward-pointing arrow from  $u^i$  to  $u^j$  indicates that  $u^i \succeq u^j$ .

1 A byproduct of Corollary 1 is that the sets  $e(\bigcup_{\pi} \{u_1^\pi\}), \dots, e(\bigcup_{\pi} \{u_N^\pi\})$  are nonempty. In  
 2 particular, there is at least one F-optimal policy.

3 **Theorem 2.** *Let  $\Pi_F^*$  be the set of all F-optimal policies. Then  $\Pi_F^* \neq \emptyset$ .*

4 This existence result can also be obtained in a different way. [Puterman \(2014\)](#) has shown  
 5 that a scalar-valued Markov decision process satisfying the assumptions of this work – namely,  
 6 a finite  $S$ , a compact  $A$ , and rewards and transition probabilities which are continuous on  $A$   
 7 – has at least one optimal policy, that is, a  $\pi^*$  such that  $u_1^{\pi^*}(s) = \max_{\pi \in \Pi} u_1^\pi(s)$  for all states  
 8  $s$ . Take now any  $m$  positive scalars  $\lambda_1, \dots, \lambda_m$ , and write  $\lambda = (\lambda_1, \dots, \lambda_m)$ . Recalling that  
 9  $R_t(s, a) = (r_t(s, a)_1, \dots, r_t(s, a)_m)$ , Puterman’s result implies that the Markov decision process  
 10 with rewards  $\langle \lambda, R_t(s, a) \rangle \in \mathbb{R}$  has a policy  $\pi^*$  such that  $\langle \lambda, u_1^{\pi^*}(s) \rangle \geq \langle \lambda, u_1^\pi(s) \rangle$  for all  $\pi \in \Pi$   
 11 and  $s \in S$ . It is a straightforward exercise to prove that such a  $\pi^*$  must be V-optimal, and  
 12 therefore F-optimal, for the vector-valued process.

1 We now propose dynamic programming equations that yield an algorithm for enumerating  
 2 the set of F-optimal policies. Our proof of the equations' validity will employ the following  
 3 observation.

4 **Lemma 3.** For all  $t = 1, \dots, N - 1$ ,

$$\bigcup_{d_t \in D} \left\{ S \ni s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))v(j) : v \in \bigcup_{\pi} \{u_{t+1}^{\pi}\} \right\} = \bigcup_{\pi} \{u_t^{\pi}\}$$

5 *Proof.* Let  $t = 1, \dots, N - 1$ . Let  $w \in F(S, \mathbb{R}^m)$  such that

$$\forall s \in S, w(s) = R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))v(j)$$

6 for some  $d_t \in D$  and  $v \in \bigcup_{\pi} \{u_{t+1}^{\pi}\}$ . We may write  $v = u_{t+1}^{\pi}$  for some  $\pi = (d_1, \dots, d_{N-1}) \in \Pi$ .  
 7 Let  $\pi' \in \Pi$  be any policy such that  $\bar{\pi}'(t) = (d_t, d_{t+1}, \dots, d_{N-1})$ . Then for all  $s \in S$ ,  $w(s) =$   
 8  $u_t^{\pi'}(s)$ . Thus,  $w = u_t^{\pi'}$ , whence

$$\bigcup_{d_t \in D} \left\{ S \ni s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))v(j) : v \in \bigcup_{\pi} \{u_{t+1}^{\pi}\} \right\} \subseteq \bigcup_{\pi} \{u_t^{\pi}\}.$$

9 The converse inclusion can readily be obtained from (3); as a result, the lemma is established.  
 10  $\square$

11 Recall that  $D$ , the set of all decision rules, is the set of all mappings from  $S$  into  $A$ . We  
 12 may now state the relation between  $e(\bigcup_{\pi \in \Pi} \{u_t^{\pi}\})$  and  $e(\bigcup_{\pi \in \Pi} \{u_{t+1}^{\pi}\})$  for all  $t = 1, \dots, N - 1$ .  
 13 From it we will deduce a dynamic programming algorithm that finds all F-optimal policies by  
 14 leveraging the structure of the equations.

15 **Theorem 3.** For all  $t = 1, \dots, N$ ,  $e(\bigcup_{\pi \in \Pi} \{u_t^{\pi}\})$  is the unique solution  $U_t$  to either of the  
 16 following equations:

$$U_t = e \left( \bigcup_{d_t \in D} \left\{ S \ni s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))v(j) : v \in U_{t+1} \right\} \right); \quad t < N \quad (9)$$

17

$$U_t = \{R_N\}; \quad t = N \quad (10)$$

18 *Proof.* For greater legibility we set, for all  $t = 1, \dots, N$ ,

$$G_t = \bigcup_{d_t \in D} \left\{ S \ni s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))v(j) : v \in U_{t+1} \right\}.$$

19 We proceed by induction on  $t$ . For any policy  $\pi \in \Pi$ ,  $u_N^{\pi} = R_N$ . Thus,  $\bigcup_{\pi} \{u_N^{\pi}\}$  is the  
 20 singleton  $\{R_N\}$ , and  $e(\bigcup_{\pi} \{u_t^{\pi}\}) = \{R_N\} = U_N$ . The property then holds for  $t = N$ . Assume  
 21 it is true for  $t + 1$ , for some  $t < N$ . Let  $u \in e(\bigcup_{\pi} \{u_t^{\pi}\})$ . There exists a  $\pi = (d_1, \dots, d_{N-1}) \in \Pi$   
 22 such that  $u = u_t^{\pi}$ . By Corollary 1, there exists a  $v \in e(\bigcup_{\pi} \{u_{t+1}^{\pi}\})$  such that  $v \succeq u_{t+1}^{\pi}$ . Let

1  $w \in F(S, \mathbb{R}^m)$  be the function such that  $w(s) = R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))v(j)$  for all  
 2  $s \in S$ . Then  $w \in \bigcup_{\pi} \{u_t^{\pi}\}$  by Lemma 3, and  $w \succeq u$  by Lemma 1. But  $u$  is efficient in  $\bigcup_{\pi} \{u_t^{\pi}\}$ ;  
 3 therefore,  $u = w$ . Appealing to the induction hypothesis for  $v$  then to Lemma 3 proves  $u \in U_t$ .  
 4 This establishes  $e(\bigcup_{\pi} \{u_t^{\pi}\}) \subseteq U_t$ .

5 To show the converse inclusion, let  $v \in U_t$ . We have  $v \in \bigcup_{\pi} \{u_t^{\pi}\}$ . Consider now some  
 6  $u \in \bigcup_{\pi} \{u_t^{\pi}\}$  such that  $u \succeq v$ . We shall show that we necessarily have  $v = u$ . Applying  
 7 Corollary 1 then Lemma 1, there is a  $u' \in e(\bigcup_{\pi} \{u_{t+1}^{\pi}\})$  such that  $w := s \mapsto R_t(s, d_t(s)) +$   
 8  $\sum_{j \in S} p_t(j|s, d_t(s))u'(j) \succeq u$ . Then  $w \succeq v$ . By our induction hypothesis,  $u' \in U_{t+1}$ , and thus  
 9  $w \in G_t$ . But  $v$  being efficient in  $G_t$ , we must have  $v = w$  and therefore  $v \succeq u$ . Consequently,  
 10  $v = u$ . The requisite inclusion then follows, and the property holds for all  $t = 1, \dots, N$ .  $\square$

11 Although Equations (9) and (10) bear a striking resemblance to the White equations (see  
 12 Section 1), the two sets of equations differ in crucial respects. In the first place, the unknowns in  
 13 (9) and (10) are subsets of  $F(S, \mathbb{R}^m)$ , whereas in White's case they are subsets of  $\mathbb{R}^m$ . White's  
 14 equations involve a total of  $N \cdot |S|$  unknowns; ours involve  $N$  unknowns. In the second place,  
 15 the solution of (9) or (10) at an epoch  $t$  must be a subset of  $\bigcup_{\pi \in \Pi} \{u_t^{\pi}\}$  by virtue of, *inter alia*,  
 16 Lemma 3, the key argument in the proof above. In contrast, we cannot guarantee in general  
 17 that White's solution sets are contained in the  $\bigcup_{\pi \in \Pi} \{u_t^{\pi}(s)\}$ 's,  $s \in S$ , except in conditions  
 18 like those expatiated in (Miframi, 2023). Incidentally, two of the conditions – namely, that  
 19 the dynamics of the model be deterministic, and that the decision making horizon be short of  
 20 three epochs – are special cases of this paper's hypotheses, and thus ensure the validity of both  
 21 White's and our equations.

22 It is clear that by construction each member of  $U_t$ ,  $t = 1, \dots, N - 1$ , is characterized by some  
 23 sequence of decision rules  $(d_t, \dots, d_{N-1})$ . Thus, if  $\mathcal{L}_t$  is the mapping from  $D^{N-t}$  into  $U_t$  defined  
 24 by  $\mathcal{L}_t(d_t, \dots, d_{N-1}) = u_t^{\pi}$ , where  $\pi$  is any policy with  $\bar{\pi}(t) = (d_t, \dots, d_{N-1})$ , then  $\mathcal{L}_t$  is onto. This  
 25 mapping need not be one-to-one, as distinct policies may well have the same return function.

26 Theorem 3 and the remarks of the previous paragraph naturally give rise to the following  
 27 algorithm for calculating F-optimal policies, an algorithm that also solves Equations (9) subject  
 28 to (10).

29 **Algorithm 1.** *Solution of Equations (9) subject to (10) and calculation of  $\Pi_F^*$ .*

30 (1) Set  $t = N - 1$  and

$$U_{N-1} = e\left(\bigcup_{d_t \in D} \left\{S \ni s \mapsto R_{N-1}(s, d_t(s)) + \sum_{j \in S} p_{N-1}(j|s, d_t(s))R_N(j)\right\}\right) \quad (11)$$

31

$$P_{N-1}^* = \{d_t \in D : S \ni s \mapsto R_{N-1}(s, d_t(s)) + \sum_{j \in S} p_{N-1}(j|s, d_t(s))R_N(j) \in U_{N-1}\} \quad (12)$$

(2) Substitute  $t - 1$  for  $t$  and set

$$U_t = e\left(\bigcup_{d_t \in D} \left\{ S \ni s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))v(j) : v \in U_{t+1} \right\}\right) \quad (13)$$

$$P_t^* = \{(d_t, d_{t+1}, \dots, d_{N-1}) \in D^{N-t} : (d_{t+1}, \dots, d_{N-1}) \in P_{t+1}^* \\ \text{and } S \ni s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))u_{t+1}^\pi(j) \in U_t\} \quad (14)$$

where  $\pi \in \Pi$  is any policy such that  $\bar{\pi}(t+1) = (d_{t+1}, \dots, d_{N-1})$ .

(3) If  $t = 1$ , stop. Otherwise, go to (2).

**Proposition 2.** The sets  $U_t$  returned by Algorithm 1 are the solutions to Equations (9) and (10), and therefore satisfy  $U_t = e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$  for each  $t = 1, \dots, N$ .

At termination,  $P_1^*$  contains every policy satisfying  $u_t^\pi \in U_t$  for all  $t = 1, \dots, N$ . By Theorem 3, such policies are F-optimal, but one might ask whether these include all or only a subset of F-optimal policies. It turns out from Lemma 1 that if  $\pi$  is F-optimal, then the subpolicy  $\bar{\pi}(t')$ ,  $1 < t' \leq N$ , is also F-optimal with respect to the portion of the decision making horizon that begins at  $t'$ , i.e  $u_{t'}^\pi \in U_{t'}$ . Phrased loosely, an F-optimal policy  $\pi$  is ‘‘F-optimal’’ at every stage of decision making: not only does it achieve an efficient return function  $u_1^\pi$  over the  $N$  epochs, but it also achieves an efficient return function  $u_t^\pi$  from any epoch  $t$  onward. Specifically, if for each epoch  $t$  we let  $E_t$  denote the efficient set  $e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$ , we have the sequence of implications

$$u_1^\pi \in E_1 \implies u_2^\pi \in E_2 \implies \dots \implies u_{N-1}^\pi \in E_{N-1},$$

or, viewed from another, logically equivalent angle,

$$u_{N-1}^\pi \notin E_{N-1} \implies u_{N-2}^\pi \notin E_{N-2} \implies \dots \implies u_1^\pi \notin E_1.$$

A formal statement and proof of this structural property follow.

**Proposition 3.** For any  $\pi \in \Pi$ ,  $\pi \in \Pi_F^*$  implies  $u_t^\pi \in e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$  for all  $t = 1, \dots, N$ .

*Proof.* Let  $\pi = (d_1, \dots, d_{N-1}) \in \Pi_F^*$ . We proceed by induction on  $t$ . By definition of  $\Pi_F^*$ ,  $u_1^\pi \in e(\bigcup_{\pi \in \Pi} \{u_1^\pi\})$ . Let  $t = 2, \dots, N - 1$  such that  $u_t^\pi \in e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$ . We will show that  $u_{t+1}^\pi \in e(\bigcup_{\pi \in \Pi} \{u_{t+1}^\pi\})$ . Assume to the contrary that  $u_{t+1}^\pi \notin e(\bigcup_{\pi \in \Pi} \{u_{t+1}^\pi\})$ . There then exists a  $\pi' \in \Pi$  such that  $u_{t+1}^{\pi'} \succ u_{t+1}^\pi$ . Therefore,  $s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))u_{t+1}^{\pi'}(j) \succ u_{t+1}^\pi$  by Lemma 1. However, Lemma 3 tells us that  $s \mapsto R_t(s, d_t(s)) + \sum_{j \in S} p_t(j|s, d_t(s))u_{t+1}^{\pi'}(j) \in \bigcup_{\pi \in \Pi} \{u_t^\pi\}$ , which contradicts the fact that  $u_t^\pi \in e(\bigcup_{\pi \in \Pi} \{u_t^\pi\})$  and completes the proof.  $\square$

The argument of this proof is very straightforward. Suppose a policy  $\pi = (d_1, \dots, d_{t+1}, \dots, d_{N-1})$  does *not* achieve an efficient return function over the period  $t + 1, \dots, N$ . This means we can

1 find a policy  $\pi'$  whose return function for the same period is better, i.e.  $u_{t+1}^{\pi'} \succ u_{t+1}^{\pi}$ . Obviously,  
 2 because decisions taken prior to epoch  $t + 1$  cannot influence how well a policy does between  
 3  $t + 1$  and  $N$ , we may write, with a slight abuse of notation,  $u_{t+1}^{\pi} = u_{t+1}^{\bar{\pi}(t+1)}$  and  $u_{t+1}^{\pi'} = u_{t+1}^{\bar{\pi}'(t+1)}$ ,  
 4 so that  $u_{t+1}^{\bar{\pi}'(t+1)} \succ u_{t+1}^{\bar{\pi}(t+1)}$ . Now suppose we extended both  $\bar{\pi}(t+1)$  and  $\bar{\pi}'(t+1)$  by the decision  
 5 rule  $d_t$ , giving rise to two (partial) policies  $(d_t, \bar{\pi}(t+1))$  and  $(d_t, \bar{\pi}'(t+1))$ . Since both policies  
 6 employ the same decision rule at time  $t$ , and since  $u_{t+1}^{\bar{\pi}'(t+1)} \succ u_{t+1}^{\bar{\pi}(t+1)}$ , we deduce from Equation  
 7 (3) that  $(d_t, \bar{\pi}'(t+1))$  has a superior return function, that is,  $u_t^{(d_t, \bar{\pi}'(t+1))} \succ u_t^{(d_t, \bar{\pi}(t+1))}$ . But  
 8  $d_t$  being the  $t$ -th decision rule in  $\pi$ , we have  $(d_t, \bar{\pi}(t+1)) = \bar{\pi}(t)$  by construction, so that  
 9  $u_t^{(d_t, \bar{\pi}(t+1))} = u_t^{\pi}$ , and thus  $\pi$  does not achieve an efficient return function over the period  $t, \dots,$   
 10  $N$ . What we have just shown, in summary, is that if  $\pi$  is not “F-optimal” from some epoch  
 11  $t + 1$  onward, it is not “F-optimal” from epoch  $t$  onward either, no matter what actions are  
 12 prescribed at epoch  $t$ .

13 **Example 4.** For each  $u^i$  in  $\bigcup_{\pi \in \Pi} \{u_1^{\pi}\}$  of Example 3, write  $u^i = u_1^{\pi_i}$ . Therefore,

$$\bigcup_{\pi \in \Pi} \{u_1^{\pi}\} = \{u_1^{\pi_1}, \dots, u_1^{\pi_8}\}.$$

14 Recall that  $e(\bigcup_{\pi \in \Pi} \{u_1^{\pi}\}) = \{u_1^{\pi_1}, u_1^{\pi_2}\}$ . The set of return functions for the period  $t = 2, \dots, N$   
 15 for this model was found to be equal to

$$\bigcup_{\pi \in \Pi} \{u_2^{\pi}\} = \{u_2^{\pi_1}, u_2^{\pi_3}, u_2^{\pi_5}, u_2^{\pi_7}\},$$

16 with  $u_2^{\pi_2} = u_2^{\pi_1}$ ,  $u_2^{\pi_4} = u_2^{\pi_3}$ ,  $u_2^{\pi_6} = u_2^{\pi_5}$ , and  $u_2^{\pi_8} = u_2^{\pi_7}$ . The relations within this set are  
 17 depicted by the diagram in Figure 2. As is obvious from this diagram,

$$e\left(\bigcup_{\pi \in \Pi} \{u_2^{\pi}\}\right) = \{u_2^{\pi_1}\} = \{u_2^{\pi_2}\},$$

18 so that  $u_1^{\pi_i} \in e(\bigcup_{\pi \in \Pi} \{u_1^{\pi}\})$  does indeed imply  $u_2^{\pi_i} \in e(\bigcup_{\pi \in \Pi} \{u_2^{\pi}\})$ . To provide further confir-  
 19 mation of Proposition 3, similar calculations were carried out for the period 3, ...,  $N$ , yielding  
 20 the result in Figure 3 that

$$e\left(\bigcup_{\pi \in \Pi} \{u_3^{\pi}\}\right) = \{u_3^{\pi_1}\} = \{u_3^{\pi_2}\} = \{u_3^{\pi_3}\} = \{u_3^{\pi_4}\}.$$

21 Again, this shows that  $u_1^{\pi_i} \in e(\bigcup_{\pi \in \Pi} \{u_1^{\pi}\})$  implies  $u_3^{\pi_i} \in e(\bigcup_{\pi \in \Pi} \{u_3^{\pi}\})$ . Observe that the  
 22 converse implication fails: for example,  $\pi_4$  achieves an efficient return function between epochs  
 23 3 and  $N$ , yet between epochs 1 and  $N$  its return function is inefficient.

24 As a result of Proposition 3, we are assured that Algorithm 1 will discover all F-optimal  
 25 policies, since the policies that satisfy  $u_1^{\pi} \in U_1 = e(\bigcup_{\pi \in \Pi} \{u_1^{\pi}\})$  are precisely those that satisfy  
 26  $u_t^{\pi} \in U_t = e(\bigcup_{\pi \in \Pi} \{u_t^{\pi}\})$  for all  $t = 1, \dots, N$ .



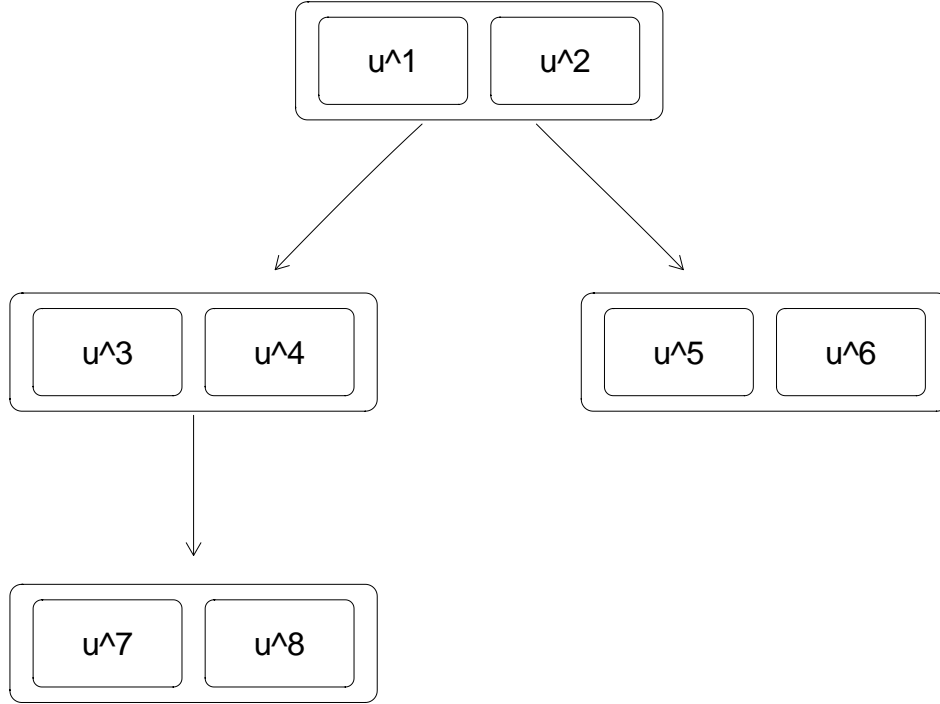


FIGURE 2. A Hasse diagram of the set  $\bigcup_{\pi \in \Pi} \{u_2^\pi\}$  in Example 4, ordered by  $\succeq$ . The arrows carry the same significance as in Figure 1. Elements occupying the same rectangle are equal. For example, the uppermost rectangle means that  $u_2^{\pi_1} = u_2^{\pi_2}$ .

1 **Corollary 2.** Algorithm 1 is guaranteed to locate all F-optimal policies at termination, i.e.  
 2  $P_1^* = \Pi_F^*$ , where  $P_t^*$  is defined by Equation (14) for each  $t = 1, \dots, N - 1$ .

3 Recall that part of the motivation for introducing F-optimality was that it is a necessary  
 4 condition for V-optimality. Hence, if  $\Pi_V^*$  denotes the set of all V-optimal policies, then  $\Pi_V^* \subseteq$   
 5  $\Pi_F^*$ . As a preliminary to an algorithm for the determination of  $\Pi_V^*$ , we show that given an  
 6 initial state  $s \in S$ , each efficient policy return vector accrued over the decision making horizon is  
 7 attained by at least an F-optimal policy. For any state  $s$ , the efficient elements in  $\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}$   
 8 are therefore a subset of  $\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\}$ . This has the practical effect of reducing the task of  
 9 “maximizing” return vectors over the whole of  $\Pi$  to the less onerous task of “maximizing”  
 10 return vectors over  $\Pi_F^*$ .

11 **Theorem 4.** Let  $s \in S$ . Then  $e\left(\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}\right) = e\left(\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\}\right)$ .

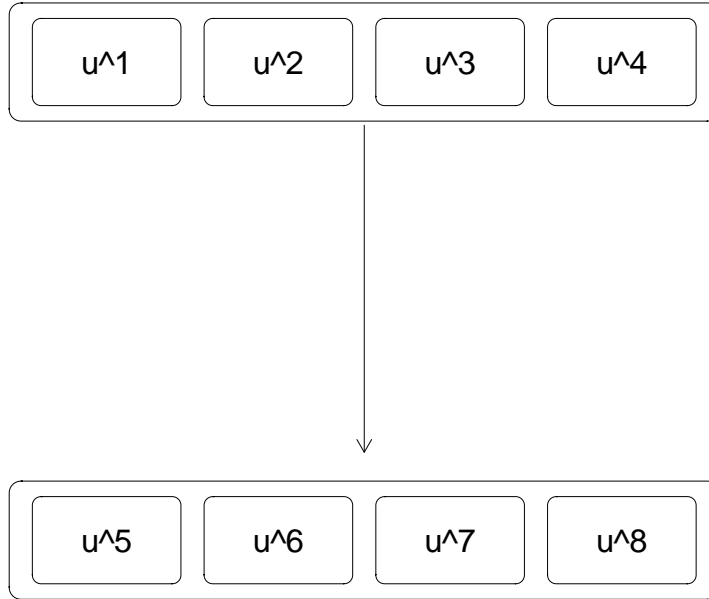


FIGURE 3. A Hasse diagram of the set  $\bigcup_{\pi \in \Pi} \{u_3^\pi\}$  in Example 4, ordered by  $\succ$ . The arrows and rectangles have the same interpretation as in Figure 2

1 *Proof.* Decompose  $\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}$  as follows:

$$\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\} = F_1 \cup F_2,$$

2 where

$$F_1 = \bigcup_{\pi \in \Pi \setminus \Pi_F^*} \{u_1^\pi(s)\}$$

3 and

$$F_2 = \bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\}.$$

4 We shall first prove that  $e\left(\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}\right) \subseteq e(F_2)$ . Pick a policy  $\pi^* \in \Pi$  such that  $u_1^{\pi^*}(s) \in$   
 5  $e\left(\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}\right)$ . It shall be established that  $u_1^{\pi^*}(s) \in F_2$ , which, given that  $F_2$  is a subset  
 6 of  $\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}$ , implies that  $u_1^{\pi^*}(s)$  is efficient in  $F_2$ . Suppose, for the sake of contradiction,  
 7 that  $u_1^{\pi^*}(s) \notin F_2$ . Then  $\pi^* \notin \Pi_F^*$ , and there is therefore, applying Theorem 1, a  $\pi' \in \Pi_F^*$  with

1  $u_1^{\pi'} \succ u_1^{\pi^*}$ . Thus  $u_1^{\pi'}(s) \geq u_1^{\pi^*}(s)$ . Since  $u_1^{\pi^*}(s)$  is efficient in  $\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}$ , it follows that  
 2  $u_1^{\pi^*}(s) = u_1^{\pi'}(s) \in F_2$ : a contradiction. Therefore,  $u_1^{\pi^*}(s) \in F_2$ , hence  $u_1^{\pi^*}(s) \in e(F_2)$ . This  
 3 concludes the demonstration of the fact that  $e\left(\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}\right) \subseteq e(F_2)$ .

4 Consider now an F-optimal policy  $\pi^* \in \Pi_F^*$  satisfying  $u_1^{\pi^*}(s) \in e(F_2)$ . To show that  $u_1^{\pi^*}(s)$   
 5 is also efficient in  $\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}$ , let  $\pi \in \Pi$  be some policy such that  $u_1^\pi(s) \geq u_1^{\pi^*}(s)$ . Either  
 6  $\pi \in \Pi_F^*$  or  $\pi \notin \Pi_F^*$ . If  $\pi \in \Pi_F^*$ , then  $u_1^\pi(s) \in F_2$  and consequently  $u_1^\pi(s) = u_1^{\pi^*}(s)$ . Otherwise,  
 7 invoking Theorem 1 again, there exists a  $\pi' \in \Pi_F^*$  such that  $u_1^{\pi'}(s) \geq u_1^\pi(s)$ . Due to the  
 8 transitivity of  $\geq$  and given that  $u_1^{\pi'}(s) \in F_2$ , this implies  $u_1^{\pi'}(s) = u_1^{\pi^*}(s)$ . It follows that  
 9  $u_1^{\pi^*}(s) \geq u_1^\pi(s)$ , which when combined with the fact that  $u_1^\pi(s) \geq u_1^{\pi^*}(s)$  yields  $u_1^\pi(s) = u_1^{\pi^*}(s)$ .  
 10 In both cases, we have that  $u_1^\pi(s) = u_1^{\pi^*}(s)$ . This proves the efficiency of  $u_1^{\pi^*}(s)$  in  $\bigcup_{\pi \in \Pi} \{u_1^\pi(s)\}$   
 11 and concludes the proof of the theorem.  $\square$

12 A characterization of V-optimal policies follows at once from this theorem, namely that a  
 13 policy  $\pi^* \in \Pi$  is V-optimal if and only if  $u_1^{\pi^*}(s) \in e\left(\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\}\right)$  for each state  $s \in S$ . We  
 14 therefore have an alternate – and, as will now be shown, useful – representation of  $\Pi_V^*$ :

$$\Pi_V^* = \left\{ \pi \in \Pi_F^* : \forall s \in S, u_1^\pi(s) \in e\left(\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\}\right) \right\}.$$

15 **Example 5.** Let us pursue Examples 3 and 4 in the light of Theorem 4. For the particular  
 16 model under consideration, it was found that

$$e\left(\bigcup_{\pi \in \Pi} \{u_1^\pi(1)\}\right) = \left\{ (1.64, 12.6), (3.84, 5.66) \right\}$$

17 and

$$e\left(\bigcup_{\pi \in \Pi} \{u_1^\pi(2)\}\right) = \left\{ (2.44, 8.56) \right\}.$$

18 It was established earlier that  $e\left(\bigcup_{\pi \in \Pi} \{u_1^\pi\}\right) = \{u_1^{\pi_1}, u_1^{\pi_2}\}$  for two policies  $\pi_1$  and  $\pi_2$ . By  
 19 definition, then,  $\pi_1$  and  $\pi_2$  are F-optimal, and for each state  $s \in S = \{1, 2\}$  we have that

$$\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\} = \{u_1^{\pi_1}(s), u_1^{\pi_2}(s)\}.$$

20 It was also found that

$$\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(1)\} = \left\{ (1.64, 12.6), (3.84, 5.66) \right\},$$

21 hence

$$e\left(\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(1)\}\right) = \left\{ (1.64, 12.6), (3.84, 5.66) \right\},$$

1 and that for  $s = 2$ ,

$$\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(2)\} = \left\{ (2.44, 8.56) \right\} = e \left( \bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(2)\} \right),$$

2 because  $u_1^{\pi_1}(2) = u_1^{\pi_2}(2) = (2.44, 8.56)$ . Thus, Theorem 4 holds.

3 Write  $S = \{s_1, \dots, s_{|S|}\}$ . Supposing Algorithm 1 was used to generate  $\Pi_F^* = \{\pi_1, \dots, \pi_n\}$ ,  
 4 which we assume here to be finite, Theorem 6 suggests and justifies the following procedure  
 5 for the calculation of V-optimal policies.

6 **Algorithm 2.** *Calculation of  $\Pi_V^*$ .*

7 *Set  $T = \Pi_F^*$ .*

8 **For**  $i \in \{1, \dots, |S|\}$  **do**:

9 **For**  $j \in \{1, \dots, n\}$  **do**:

10 **For**  $k \in \{1, \dots, n\} \setminus \{j\}$  **do**:

11 **If**  $u_1^{\pi_k}(s_i) > u_1^{\pi_j}(s_i)$ , **drop**  $\pi_j$  from  $T$ .

12 **Proposition 4.** *Algorithm 2 terminates with  $T = \Pi_V^*$ .*

13 *Proof.* At termination, a policy  $\pi$  is in  $T$  if and only if  $\forall s \in S, \forall \pi' \in \Pi_F^*, u_1^{\pi'}(s) \not> u_1^\pi(s)$ .  
 14 Thus,  $\pi \in T$  if and only if  $\forall s \in S, u_1^\pi(s) \in e(\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\})$ . Consequently,  $T = \Pi_V^*$  by the  
 15 latter set's alternate representation.  $\square$

#### 16 4. NUMERICAL EXPERIMENTS

17 The algorithms that have been developed were tested on randomly-generated instances of  
 18 a three-state, two-action model with six epochs and varying numbers of criteria ( $m$ ). In each  
 19 instance, the rewards and transition probabilities associated with an epoch  $t$  were sampled from  
 20 exponential distributions then, in the case of the probabilities, scaled to  $[0, 1]$ .

21 Algorithm 1 was used to locate the efficient return functions in  $\bigcup_{\pi \in \Pi} \{u_1^\pi\}$ , in addition to  
 22 the corresponding F-optimal policies, through solving Equations (9) and (10) for  $U_N = U_6, U_5,$   
 23  $U_4, U_3, U_2$  then  $U_1$ . Algorithm 2 was used afterwards to determine which of the F-optimal  
 24 policies were V-optimal.

25 To implement steps (1) and (2) of Algorithm 1, full search was used given the finite number  
 26 of decision rules available. The experiments were programmed in C and run on a quadcore  
 27 Intel Core i5-1145G7 laptop with 16GB of RAM.

28 The results are collected in Table 1, where the experiments are grouped by  $m$  into ten  
 29 groups of a hundred experiments each. Bearing in mind that  $U_1 = e(\bigcup_{\pi \in \Pi} \{u_1^\pi\})$  (Proposition  
 30 2), the column  $|U_1|$  reports the minimum and maximum value observed in a single group of  
 31 experiments of the number of efficient points in  $\bigcup_{\pi \in \Pi} \{u_1^\pi\}$ . Similarly, the column  $|\Pi_F^*|$  contains  
 32 the range of values taken by the number of F-optimal policies for a single group. Immediately

1 to the right is a column indicating the range of CPU time expended on Algorithm 1. The last  
 2 two columns provide analogous statistics for Algorithm 2, with  $|\Pi_V^*|$  indicating the minimum  
 3 and maximum number of V-optimal policies identified in a group.

$m$	$ U_1 $	$ \Pi_F^* $	Algorithm 1 (seconds)	$ \Pi_V^* $	Algorithm 2 (seconds)
1	1-2	1-2	0	1	0
2	11-253	11-253	0-0.0280	1-85	0
3	217-3903	217-3903	0.0120-2.4238	70-874	0-0.0040
4	853-7580	853-7580	0.2199-5.8955	249-6496	0
5	1259-7159	1259-7159	0.1440-6.7715	790-3345	0-0.0040
6	3817-21701	3817-21701	1.5478-39.5032	1726-18190	0-0.0040
7	2230-18169	2230-18169	0.4879-27.0269	1174-15611	0-0.0040
8	5668-21874	5668-21874	3.5077-56.9811	5123-14579	0-0.0040
9	2972-25937	2972-25937	1.8531-60.5031	2395-25937	0-0.0080
10	8005-27636	8005-27636	10.0394-93.4665	6558-27636	0-0.0079

TABLE 1. Results for randomly-generated problems grouped by  $m$ . The data is presented in minimum-maximum format.

4 As one would have predicted, the sizes of  $U_1$ ,  $\Pi_F^*$  and  $\Pi_V^* \subseteq \Pi_F^*$  tended to grow as  $m$   
 5 increased. The number of F-optimal policies ranged from one policy for a scalar-valued problem  
 6 to 27 636 policies for a problem with ten objectives. The fraction of F-optimal policies that  
 7 were V-optimal rose together with  $m$ , although V-optimal policies accounted in the majority  
 8 of experiments for less than half of  $\Pi_F^*$ .

9 The computational demands for solving Equations (9) and (10) then finding every policy in  
 10  $\Pi_F^*$  and  $\Pi_V^*$  also tended to grow with  $m$ . For  $m \leq 4$ , CPU time ranged from zero to six seconds  
 11 for the two algorithms combined. Given that a typical multi-objective decision making problem  
 12 seldom exceeds four objectives (Stewart, Palmer, & DuPont, 2021), these results suggest that  
 13 the algorithms possess the potential for being effective solution methods in a wide array of real  
 14 world applications.

15 It should be noted that, relative to Algorithm 1, Algorithm 2 required negligible amounts  
 16 of time due to the alternate representation of  $\Pi_V^*$  arrived at in the previous section. Recall  
 17 that a policy  $\pi$  is V-optimal by definition if for all states  $s$  the vector  $u_1^\pi(s)$  is efficient in  
 18  $\bigcup_\pi \{u_1^\pi(s)\}$ . This is equivalent, as we have seen, to efficiency of each  $u_1^\pi(s)$  merely in the  
 19 subset  $\bigcup_{\pi \in \Pi_F^*} \{u_1^\pi(s)\}$ . To fully appreciate the practicality of this alternate representation, we  
 20 measured, during each of the previous experiments, the time needed to locate all V-optimal  
 21 policies through an exhaustive search of the sets  $\bigcup_\pi \{u_1^\pi(s)\}$ ,  $s \in S$ . In Table 2 we juxtapose  
 22 the results against those already reported for Algorithm 2.

$m$	Exhaustive search (seconds)	Algorithm 2 (seconds)
1	10.4838-11.0478	0
2	15.5674-30.9460	0
3	26.7423-39.8178	0
4	31.0000-53.0195	0
5	54.9910-84.9934	0-0.0040
6	50.0692-98.9275	0-0.0040
7	79.7329-121.6538	0-0.0040
8	91.9299-146.0090	0-0.0040
9	94.4560-126.2161	0-0.0080
10	148.1788-186.5719	0-0.0079

TABLE 2. Comparison of the execution times of two methods of enumerating the set of V-optimal policies: exhaustive search of  $\Pi$  and Algorithm 2. The problems considered here are those summarized in Table 1.

1 It is clear from the table that Algorithm 2, which relies on the alternate representation of  
2  $\Pi_V^*$ , was consistently and substantially faster than a strategy based on direct optimization over  
3  $\Pi$ . The gap between the two methods became more pronounced as  $m$  increased. More impor-  
4 tantly, whereas the time required by exhaustive search multiplied more than tenfold between  
5  $m = 1$  and  $m = 10$ , Algorithm 2 ran at a comparatively stable speed.

6

7

## 5. APPLICATION TO INVENTORY CONTROL

8 As a further illustration of our theoretical results, we consider the stochastic inventory  
9 control problem described in (Puterman, 2014, p. 38). The problem will be restated here  
10 for completeness. We take the position of a warehouse manager overseeing the inventory on  
11 hand of a particular product. Based on the inventory level at the beginning of each month,  
12 the manager may elect to order additional stock so long as it does not exceed the warehouse's  
13 capacity,  $M$ . The manager must keep sufficient inventory to meet external demand, but must  
14 also avoid overordering stock so as to minimize storage and ordering costs.

15 During each month, the events unfold as follows: (1) the decision is made to purchase (or  
16 not) additional stock; (2) the order is instantly fulfilled; (3) customer demand for the product  
17 arrives; then (4) if inventory is sufficient, all the demand is met on the last day of the month.  
18 External demand in month  $t$ ,  $D_t$ , follows a time-homogeneous distribution  $p_j = P(D_t = j)$ ,  
19 for all non-negative integers  $j$ . The cost of ordering  $u$  units in any month is  $O(u) = K + c(u)$   
20 if  $u > 0$  and 0 if  $u = 0$ , where  $K > 0$  is a fixed cost and  $c$  a nondecreasing function of  $u$ . The  
21 cost of holding  $u$  units between delivery of additional stock and sale of inventory at the end

1 of a month is  $h(u)$ ,  $h$  being a nondecreasing function of  $u$ . If the demand is for  $j$  units and  
 2 sufficient inventory is available, then the revenue from selling those  $j$  units is  $f(j)$ , where  $f$  is  
 3 nondecreasing in  $j$ .

4 Puterman proposes the following Markov decision process formulation. States represent the  
 5 inventory level on the first day of a month, and actions represent the amount of stock the  
 6 manager can order each month. In our notation,

$$S = \{0, \dots, M\} \quad (15)$$

7 and

$$A_s = \{0, \dots, M - s\} \quad (16)$$

8 for all  $s \in S$ . Thus, for example, if there are 3 units in the inventory for a warehouse capacity  
 9 of 5, the manager can order up to two units of stock.

10 Let  $s_t$  be the state of the inventory in month  $t$ , and  $a_t$  the amount of stock ordered in that  
 11 month. Clearly,  $s_{t+1}$  depends only on  $s_t$ ,  $a_t$  and the random demand  $D_t$ . Then for any month  
 12  $t$  and state  $j \in S$ :

$$p_t(j|s_t, a_t) = \begin{cases} 0 & \text{if } M \geq j > s_t + a_t \\ p_{s_t+a_t-j} & \text{if } M \geq s_t + a_t \geq j > 0 \\ q_{s_t+a_t} & \text{if } M \geq s_t + a_t \text{ and } j = 0 \end{cases} \quad (17)$$

13 where  $q_{s_t+a_t} = 1 - \sum_{k=0}^{s_t+a_t-1} p_k$ . Details of the derivation of (17) are supplied in (Puterman,  
 14 2014, p. 40).

15 The manager's objective in the original formulation is to maximize the difference between the  
 16 expected revenue made over  $N$  months and the expected holding and ordering costs incurred in  
 17 that period. Consequently, the reward received in month  $t$  is taken to be  $F(s_t + a_t) - O(a_t) -$   
 18  $h(s_t + a_t)$ , where  $F(u)$  is the expected monthly revenue when the stock level prior to demand  
 19 is  $u$ . Details of the derivation of  $F$  from  $f$  and the  $p_j$ 's are also provided in (Puterman, 2014,  
 20 p. 39). We assume that the value of the inventory at the start of month  $N$  is nil.

21 Our only difference with the formulation above is that we treat revenue and costs as com-  
 22 peting optimization objectives. Thus, we define the reward accrued in month  $t$  as:

$$R_t(s_t, a_t) = (F(s_t + a_t), -O(a_t) - h(s_t + a_t)) \quad (18)$$

23 with  $R_N(s_N) = (0, 0)$ . States, actions and transition probabilities are unaffected.

24 For comparative purposes, we take the same parameter values as (Puterman, 2014, p. 38):  
 25  $N = 4$  months,  $M = 3$ ,  $K = 4$ ,  $c(u) = 2u$ ,  $h(u) = u$ , and  $f(u) = 8u$ . Demand is distributed  
 26 according to  $p_0 = \frac{1}{4}$ ,  $p_1 = \frac{1}{2}$ ,  $p_3 = \frac{1}{4}$  and  $\forall j > 3, p_j = 0$ . Revenue is 8 per unit sold, inventory  
 27 holding cost is 1 per unit, placing an order costs 4, and additional stock costs 2 per unit. The  
 28 warehouse's capacity is 3 units at a time, and the planning horizon is of 4 months.

1 There is a total of  $(M + 1)! = 24$  decision rules in this model. Over three decision epochs,  
 2 this gives rise to 13824 policies. Algorithm 1 returned 1506 F-optimal policies, one of which is  
 3 given in Table 3. Column one represents the inventory level at the beginning of a month,  $s$ .  
 4 Columns two through four represent, respectively,  $d_1(s)$ ,  $d_2(s)$  and  $d_3(s)$ . Column five reports  
 5 the value of  $u_1^\pi(s)$ . For example, if we start the first month with one item in stock, this policy  
 6 recommends that two items be ordered at the end of the first and second months, then no item  
 7 be ordered at the end of the third month. This would guarantee, in expectation, a revenue of  
 8 16 against a total cost of 12.7.

9 A subsequent comprehensive search of the policy space confirmed equality between the set  
 10 of F-optimal policies,  $\Pi_F^*$ , and the set generated by Algorithm 1. This fact agrees with the  
 11 conclusion that was drawn in Corollary 2 as to the ability of Algorithm 1 to identify all F-  
 12 optimal policies.

TABLE 3. Example of an F-optimal policy  $\pi = (d_1, d_2, d_3)$  for the stochastic inventory problem.

Start of month inventory	Order 1	Order 2	Order 3	(Exp. Revenue, - Exp. Costs)
0	1	2	0	(16.0, -14.7)
1	2	2	0	(16.0, -12.7)
2	0	1	0	(16.0, -6.7)
3	0	0	0	(22.0, -11.9)

13 As a test of the validity of Algorithm 2, we compared its output,  $T = S_{\Pi_V^*}$ , with the set  
 14 of all V-optimal policies,  $\Pi_V^*$ , obtained after a full search. Of the 1506 F-optimal policies  
 15 mentioned earlier, 61 were V-optimal. The 61 matched the policies in  $T$ , thus corroborating  
 16 the characterization of  $\Pi_V^*$  provided at the end of the previous section. Unsurprisingly, the  
 17 “never order” strategy, which incurs minimal costs over the 4 months, was V-optimal.

18 It is instructive to compare the V-optimal policies constructed by solving the problem in  
 19 its original scalar-valued formulation (Puterman, 2014, p. 96) with those obtained here. In  
 20 fact, Puterman solves his problem under the same set of parameters  $(N, M, K, f, g, h, c)$   
 21 as ours, and concludes that the unique V-optimal policy is a nonstationary  $(\Sigma, \sigma)$  strategy,  
 22 namely a policy of the form: “if units in the inventory are below  $\sigma$  at the start of the month,  
 23 order enough stock to reach  $\Sigma$  units; otherwise, do not order” (Puterman, 2014, p. 38). We  
 24 make three comments in this connection. First, whereas the optimal policy is unique in the  
 25 scalar setting, there are 61 such policies in the vector setting. Second,  $\Pi_V^*$  contains both  $(\Sigma, \sigma)$   
 26 and non- $(\Sigma, \sigma)$  policies, only one of which is stationary (the “never order” policy). Third,  
 27 Puterman’s optimal policy, which was selected to maximize the difference between revenue and  
 28 cost, is also V-optimal. To see this, recall our comment in the discussion following Theorem 2



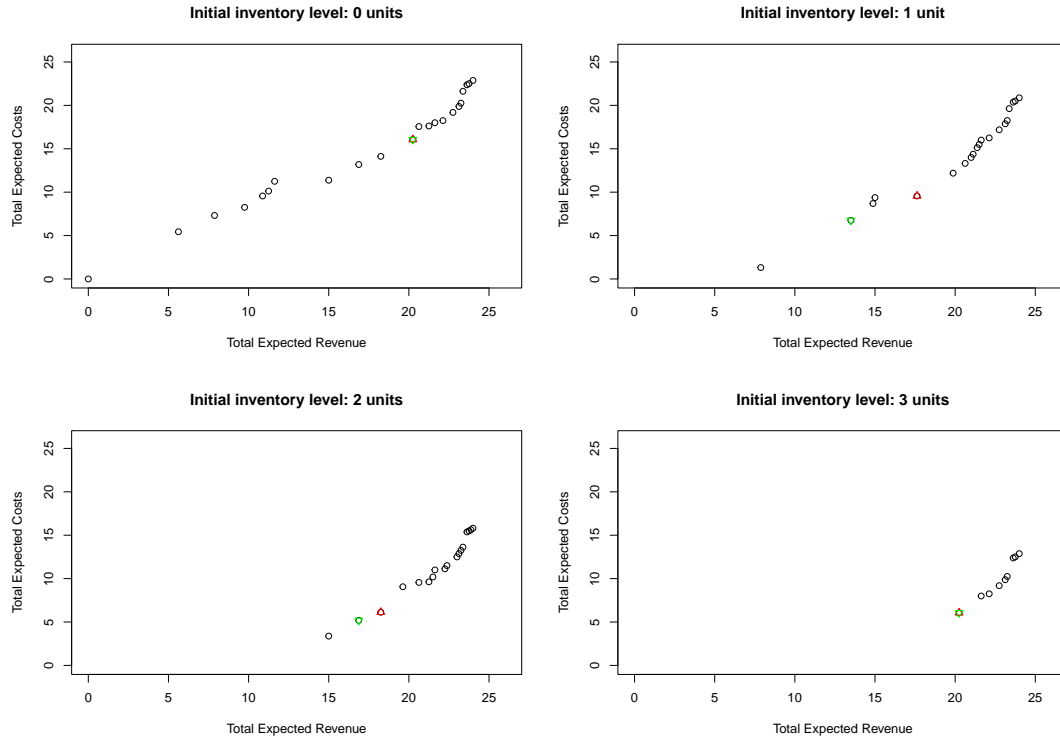


FIGURE 4. Graphical depiction of  $\bigcup_{\pi \in \Pi_V^*} \{u_1^\pi(s)\}$ ,  $s = 0, 1, 2, 3$ . A point represents one or several V-optimal policies. Sixty-one policies are represented in each plot. Table 1's policy is highlighted in green, while Puterman's optimal  $(\Sigma, \sigma)$  policy is indicated in red.

1 that a policy optimal with respect to a positive linear combination of the reward components,  
 2 i.e a  $\pi^*$  satisfying  $\sum_{i=1}^m \lambda_i u_1^{\pi^*}(s)_i \geq \sum_{i=1}^m \lambda_i u_1^\pi(s)_i$  for some positive weights  $\lambda_1, \dots, \lambda_m > 0$ ,  
 3 is V-optimal for the vector-valued model. In Puterman's case,  $\lambda_1 = \lambda_2 = 1 > 0$ , and the  
 4 conclusion follows (note that the cost component of our rewards is given by the negative of  
 5 the physical cost, so that the sum of the two components equals, in real terms, the difference  
 6 between revenue and cost).

7 For each stock level  $s \in S$ , Figure 4 portrays the returns achieved by the V-optimal policies  
 8 over the  $N = 4$  months if the initial inventory level is  $s$  units; that is, each plot depicts  
 9  $\bigcup_{\pi \in \Pi_V^*} \{u_1^\pi(s)\}$  for some  $s \in S$ . A point corresponds to one or several V-optimal policies. It  
 10 should be stressed here that  $\bigcup_{\pi \in \Pi_V^*} \{u_1^\pi(s)\}$ ,  $s \in S$ , is an efficient set by construction. This  
 11 means that given two V-optimal policies  $\pi_1$  and  $\pi_2$  yielding different returns from an inventory  
 12 level  $s$ , either  $\pi_1$  generates a higher revenue while incurring greater costs than  $\pi_2$ , or vice versa.

1 Therefore, each policy represented in  $\bigcup_{\pi \in \Pi_V^*} \{u_1^\pi(s)\}$  expresses a particular tradeoff between  
 2 costs and revenue. This is clearly reflected in all four plots.

3 Some final remarks on such plots as those of Figure 4 seem in order. When a decision is to  
 4 be made as to what policy should be enacted among those supplied by such plots, the fact that  
 5  $\bigcup_{\pi \in \Pi_V^*} \{u_1^\pi(s)\}$  is an efficient set means the decision maker has some latitude. For example, if  
 6 among the  $m$  objectives there is a high priority objective, the decision maker will prefer policies  
 7 that realize most gains in that objective. If, for example, the system has to operate under con-  
 8 straints, the decision maker will, when feasible, discard policies that violate these constraints.  
 9 Moreover, access to such plots allows the decision maker to locate policies where small con-  
 10 cessions in one objective produce considerable improvements in others. Thus, for instance, a  
 11 car manufacturer may learn that a slight increase in costs allows for substantial reduction in  
 12 tailpipe emissions, or a call center may realize that throughput may be greatly enhanced by  
 13 hiring one additional employee. Examples such as these suggest that better informed decisions  
 14 can be made when plots of  $\bigcup_{\pi \in \Pi_V^*} \{u_1^\pi(s)\}$  are available.

15

16

## 6. CONCLUSIONS AND DISCUSSION

17 To summarize, this paper endeavored to solve a class of vector-valued Markov decision pro-  
 18 cesses within two frameworks: (1) a policy is V-optimal if it delivers a maximal return from any  
 19 initial state; and (2) a policy is F-optimal if its return function over the total decision making  
 20 horizon is maximal among all return functions. An exact dynamic programming algorithm  
 21 was proposed for the second framework, which helped provide the basis for a procedure for  
 22 calculating all V-optimal policies. Fundamental to the procedure were, first, the insight that  
 23 framework (1) is subsumed under (2), and second, that a computationally useful representation  
 24 of the set of V-optimal policies can be derived from this connection. Investigation of the set  
 25 of F-optimal policies revealed that it satisfies a certain property which ensures the discovery  
 26 of all such policies by the exact algorithm. The algorithms were illustrated with numerical  
 27 experiments and a bi-objective variant of a stochastic inventory management problem.

28 In an effort to simplify the exposition, we have restricted ourselves to models with additive  
 29 rewards, but our results extend to the multiplicative case provided that the rewards meet  
 30 further assumptions. Specifically, let  $x \circ y = (x_i y_i)_{1 \leq i \leq m}$  denote the componentwise product of  
 31  $x$  and  $y$  for any  $x, y \in \mathbb{R}^m$ , and let  $u_t^\pi(s) = \mathbb{E}_\pi^s [R_t(X_t, d_t(X_t)) \circ \dots \circ R_{N-1}(X_{N-1}, d_{N-1}(X_{N-1})) \circ$   
 32  $R_N(X_N)]$  be the expected total reward for using  $\pi = (d_1, \dots, d_{N-1}) \in \Pi$  from  $t$  onward assuming  
 33 the state at this epoch is  $s$ . Supposing then that  $R_t(s, a)$  (resp.,  $R_N(s)$ ) has only nonnegative  
 34 components for all  $s \in S$ ,  $a \in A$  and  $t = 1, \dots, N-1$  (resp., for all  $s \in S$ ), every proposition in  
 35 Section 3, except Lemma 3 and Theorem 3, follows without changes. A formal proof does not  
 36 seem befitting at this stage of the paper, but the key points are these: (1) it is straightforward to  
 37 check that  $u_t^\pi(s) = R_t(s, d_t(s)) \circ \sum_{j \in S} p_t(j|s, d_t(s)) u_{t+1}^\pi(j)$  for all  $s \in S$ ,  $\pi = (d_1, \dots, d_{N-1}) \in \Pi$ ,

1  $t = 1, \dots, N - 1$ ; (2) expanding the sum in the previous expression yields the analogue of  
 2 Equation (5) where “ $\circ$ ” is substituted for “ $+$ ”; and (3) with nonnegative reward components,  
 3  $\circ$  preserves inequalities with respect to  $\geq$ . Points (1) and (3) suffice to prove Lemma 1. The  
 4 resulting expression in point (2) suffices to show, as in the original proof of Theorem 1, that  $f_t$   
 5 is upper semicontinuous for all  $t \leq N$ , and therefore that Theorem 1 is correct. Proposition 1,  
 6 Lemma 2, Birkhoff’s theorem, Theorem 4 and Proposition 4 are unrelated to whether rewards  
 7 are multiplicative or additive, and thus follow independently. The alternate representation of  
 8  $\Pi_V^*$  follows from Proposition 1 and Theorem 4. Corollary 1 is a consequence of Theorem 1.  
 9 Theorem 2 follows from Corollary 1. The analogue of Lemma 3 where “ $\circ$ ” replaces “ $+$ ” on  
 10 the left-hand side of the equality can be proven without difficulty, using the same arguments  
 11 as the original proof. From this follows Proposition 3. Theorem 3 follows from the analogue  
 12 of Lemma 3 of Theorem 1, with Equation (9) modified to reflect the change in Lemma 3. As  
 13 Algorithm 1 and Proposition 2 are based entirely on Theorem 3, they remain valid *mutatis*  
 14 *mutandis*. From the validity of Algorithm 1 follows that of Algorithm 2, and Proposition 4  
 15 holds. Finally, that Algorithm 1 is capable of finding all F-optimal policies, a result stated in  
 16 Corollary 2, is an immediate consequence of Proposition 3.

17 A difficulty to be encountered when implementing Algorithm 1 is in the computation of  
 18 efficient sets. In the numerical experiments as well as in the inventory problem, we used  
 19 enumeration because the number of actions, and therefore the number of decision rules, was  
 20 finite. If there is an infinity of actions, then analytic methods for determining the efficient  
 21 return functions in Algorithm 1 may be required. As regards Algorithm 2, the requirement  
 22 that  $\Pi_F^*$  be finite may be fulfilled even under a finite action space. However, an infinite  $\Pi_F^*$   
 23 would mean that an analytic alternative to Algorithm 2 would be in order. These questions  
 24 will be the object of a future paper.

25 Finally, it remains to consider potential applications of these results in areas beyond inven-  
 26 tory management. The operations research literature is replete with decision making problems  
 27 that can be treated as instantiations of the model studied here. For example, [Chanson, Puter-](#)  
 28 [man, and Wong \(1989\)](#) look at the problem of controlling the number of jobs that are processed  
 29 at any given moment by a computer system. Allowing too many jobs in memory can cause ex-  
 30 cessive competition for resources and hence considerable deterioration of system performance.  
 31 Just how many and *which* jobs are admitted to memory, where execution occurs, is determined  
 32 by what the authors call a “load control” policy. They assume two classes of jobs, batch and  
 33 interactive, and seek policies which minimize a weighted sum of the number of batch jobs and  
 34 the number of interactive jobs in the system. They formulate the problem as a Markov decision  
 35 process, taking as their concept of state the number of jobs in each class present in the system  
 36 as well as the fraction of those jobs occupying memory. The actions permissible at any given  
 37 time for any given state are either to admit batch jobs or interactive jobs to memory. The

1 transition probabilities are derived from a simple queueing model, and the instantaneous cost  
2 of an action is defined as a weighted sum of the current number of jobs in each class. Assigning  
3 a greater weight to either class places greater emphasis on it during minimization.

4 Chanson, Puterman and Wong’s approach to this essentially multi-objective problem is  
5 sometimes called the weighting factor approach. A different approach, which we might call the  
6 vector maximization approach, would be to change their formulation slightly, defining costs  
7 not as weighted combinations of the two classes of jobs but as unweighted two-dimensional  
8 vectors (one component per class). In our notation, shifting the point of view from costs to  
9 rewards, this would translate to  $R_t(s, a) = (-N_1(t), -N_2(t))$  for each  $(t, s, a)$ , where  $N_i(t)$   
10 equals the number of class  $i = 1, 2$  jobs in the system at time  $t$ . Optimal load control policies  
11 could then be calculated by the algorithms presented here and submitted to a decision maker  
12 for consideration. To aid the decision maker in selecting an appropriate policy, it would be  
13 advisable to supply them with plots like those of Figure 4.

14 Similar problems abound elsewhere. In medicine, for example, [Denton, Kurt, Shah, Bryant,](#)  
15 [and Smith \(2009\)](#) study the question whether and when to begin administering statins (drugs)  
16 for the treatment of lipid abnormalities in diabetes patients. The promise of increased life  
17 expectancy resulting from statin therapy encourages the initiation of such therapy at an early  
18 age. On the other hand, the sheer cost of treatment may discourage taking statins at an  
19 early age. In choosing the statin start time, therefore, a tradeoff is to be found between two  
20 competing criteria: “[the] expected future quality-adjusted time...and the annual cost of statin  
21 treatment and the cost associated with the treatment of cardiovascular events” ([Denton et al.,](#)  
22 [2009](#), p. 2). In light of these aims and of the probabilistic aspects of the problem, a Markov  
23 decision process is used.

24 The decision whether to initiate treatment is based on such changing risk factors as blood  
25 pressure, cholesterol and high-density lipoprotein. Those factors define the state of the patient  
26 at any given time. The patient is observed at periodic, discrete times at which the decision  
27 maker may elect to initiate statins or delay treatment by one period. Transitions between states  
28 depend on the particular states, the stage of treatment, and whether the patient is on statins.  
29 A longitudinal medical record obtained from a clinic treating diabetes was used in conjunction  
30 with cardiovascular risk models to estimate the probabilities of these transitions. Given a  
31 patient, the study seeks policies which maximize the expected long-term quality-adjusted life  
32 years minus therapy costs over the patient’s future. The rewards are expressed accordingly as  
33 the difference between a current estimate of quality-adjusted life years (in monetary value),  
34 which is a function of the state, and a cost component, which is a function both of the state  
35 and whether therapy has been initiated. Costs are broken down into the cost of statins and the  
36 costs associated with the treatment and follow-up of patients after a stroke or coronary heart  
37 disease event.

1 An alternative formulation – one which stems naturally from the authors’ premise that the  
2 two criteria are “competing” – would be to treat quality-adjusted life years and cost as two  
3 components of a vector-valued reward function. We might also refine this formulation by sep-  
4 arating statin costs from stroke- or coronary heart disease-related costs, thus obtaining three  
5 (or four) rather than two concurrent objectives. Under either formulation, we have finite state  
6 and action spaces and a finite horizon, so that the basic assumptions of this work are met. This  
7 means the algorithms could be employed to recommend for each patient a series of options for  
8 when to start taking statins. The physician would compare these options then select the one  
9 that is most representative of the tradeoffs they are prepared to make between the criteria.

10

## 11 DATA AVAILABILITY STATEMENT

12 Data sharing is not applicable to this article as no new data were created or analyzed in this  
13 study.

14

## 15 DISCLOSURE OF INTEREST

16 The authors have no competing interests to disclose.

17

## 18 REFERENCES

- 19 Bellman, R. (1954). The theory of dynamic programming. *Bulletin of the American Mathe-*  
20 *matical Society*, 60(6), 503–515.
- 21 Berge, C. (1985). *Graphs and Hypergraphs*. Elsevier Science Ltd.
- 22 Birkhoff, G. (1940). *Lattice Theory* (Vol. 25). American Mathematical Soc.
- 23 Borrero, J., & Akhavan-Tabatabaei, R. (2013). Time and inventory dependent optimal main-  
24 tenance policies for single machine workstations: An MDP approach. *European Journal*  
25 *of Operational Research*, 228(3), 545–555.
- 26 Brown, T. A., & Strauch, R. E. (1965). Dynamic programming in multiplicative lattices.  
27 *Journal of Mathematical Analysis and Applications*, 12(2), 364–370.
- 28 Burns, L. D., Hall, R. W., Blumenfeld, D. E., & Daganzo, C. F. (1985). Distribution strategies  
29 that minimize transportation and inventory costs. *Operations Research*, 33(3), 469–490.
- 30 Chanson, S. T., Puterman, M. L., & Wong, W. C. (1989). A Markov decision process model for  
31 computer system load control. *INFOR: Information Systems and Operational Research*,  
32 27(3), 387–402.
- 33 Chen, D. Z., Trevizan, F., & Thiébaux, S. (2023). Heuristic search for multi-objective prob-  
34 abilistic planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*  
35 (Vol. 37, pp. 11945–11954).
- 36 Coldman, A. J., & Murray, J. (2000). Optimal control for a stochastic model of cancer  
37 chemotherapy. *Mathematical Biosciences*, 168(2), 187–200.

- 1 Denton, B. T., Kurt, M., Shah, N. D., Bryant, S. C., & Smith, S. A. (2009). Optimizing the  
 2 start time of statin therapy for patients with diabetes. *Medical Decision Making*, *29*(3),  
 3 351–367.
- 4 Furukawa, N. (1980). Characterization of optimal policies in vector-valued Markovian decision  
 5 processes. *Mathematics of Operations Research*, *5*(2), 271–279.
- 6 Geoffrion, A. M. (1968). Proper efficiency and the theory of vector maximization. *Journal of*  
 7 *Mathematical Analysis and Applications*, *22*(3), 618–630.
- 8 Goedhart, J., Haijema, R., Akkerman, R., & de Leeuw, S. (2023). Replenishment and fulfilment  
 9 decisions for stores in an omni-channel retail network. *European Journal of Operational*  
 10 *Research*.
- 11 Golabi, K., Kulkarni, R. B., & Way, G. B. (1982). A statewide pavement management system.  
 12 *Interfaces*, *12*(6), 5-21.
- 13 Hayes, C. F., Rădulescu, R., Bargiacchi, E., Källström, J., Macfarlane, M., Reymond, M., ...  
 14 others (2022). A practical guide to multi-objective reinforcement learning and planning.  
 15 *Autonomous Agents and Multi-Agent Systems*, *36*(1), 26.
- 16 Henig, M. I. (1983). Vector-valued dynamic programming. *SIAM Journal on Control and*  
 17 *Optimization*, *21*(3), 490–499.
- 18 Henig, M. I. (1985). The Principle of Optimality in dynamic programming with returns in  
 19 partially ordered sets. *Mathematics of Operations Research*, *10*(3), 462–470.
- 20 Mandow, L., Pérez-de-la Cruz, J.-L., & Pozas, N. (2022). Multi-objective dynamic program-  
 21 ming with limited precision. *Journal of Global Optimization*, *82*(3), 595–614.
- 22 Mason, J., Denton, B., Shah, N., & Smith, S. (2014). Optimizing the simultaneous management  
 23 of blood pressure and cholesterol for type 2 diabetes patients. *European Journal of*  
 24 *Operational Research*, *233*(3), 727–738.
- 25 Mifrani, A. (2023). A counterexample and a corrective to the vector extension of the Bellman  
 26 equations of a Markov decision process. *arXiv preprint arXiv:2306.16937*.
- 27 Morin, T. L. (1982). Monotonicity and the Principle of Optimality. *Journal of Mathematical*  
 28 *Analysis and Applications*, *88*(2), 665–674.
- 29 Puterman, M. L. (2014). *Markov Decision Processes: Discrete Stochastic Dynamic Program-*  
 30 *ming*. John Wiley & Sons.
- 31 Ramirez-Nafarrate, A., Hafizoglu, A. B., Gel, E. S., & Fowler, J. W. (2014). Optimal control  
 32 policies for ambulance diversion. *European Journal of Operational Research*, *236*(1),  
 33 298–312.
- 34 Roijers, D., Röpke, W., Nowe, A., & Radulescu, R. (2021, July 14). On following Pareto-  
 35 optimal policies in multi-objective planning and reinforcement learning.. Retrieved from  
 36 <http://modem2021.cs.nuigalway.ie/> (Multi-Objective Decision Making Workshop  
 37 2021, MODEM 2021 ; Conference date: 14-07-2021 Through 16-07-2021)
- 38 Ruiz-Montiel, M., Mandow, L., & Pérez-de-la Cruz, J.-L. (2017). A temporal difference method

- 1 for multi-objective reinforcement learning. *Neurocomputing*, 263, 15–25.
- 2 Schlosser, R., & Gönsch, J. (2023). Risk-averse dynamic pricing using mean-semivariance  
3 optimization. *European Journal of Operational Research*, 310(3), 1151–1163.
- 4 Stewart, R. H., Palmer, T. S., & DuPont, B. (2021). A survey of multi-objective optimization  
5 methods and their applications for nuclear scientists and engineers. *Progress in Nuclear  
6 Energy*, 138, 103830.
- 7 Van Moffaert, K., & Nowé, A. (2014). Multi-objective reinforcement learning using sets of  
8 Pareto dominating policies. *The Journal of Machine Learning Research*, 15(1), 3483–  
9 3512.
- 10 Wang, L., Demeulemeester, E., Vansteenkiste, N., & Rademakers, F. E. (2024). Capacity  
11 and surgery partitioning: An approach for improving surgery scheduling in the inpatient  
12 surgical department. *European Journal of Operational Research*, 313(1), 112–128.
- 13 White, D. J. (1982). Multi-objective infinite-horizon discounted Markov decision processes.  
14 *Journal of Mathematical Analysis and Applications*, 89(2), 639–647.
- 15 White, D. J. (1993). A survey of applications of Markov decision processes. *Journal of the  
16 Operational Research Society*, 44(11), 1073–1096.
- 17 Wiering, M. A., & De Jong, E. D. (2007). Computing optimal stationary policies for multi-  
18 objective Markov decision processes. In *2007 IEEE International Symposium on Approx-  
19 imate Dynamic Programming and Reinforcement Learning* (pp. 158–165).
- 20 Zadeh, L. (1963). Optimality and non-scalar-valued performance criteria. *IEEE Transactions  
21 on Automatic Control*, 8(1), 59–60.
- 22 Zorn, M. (1935). A remark on method in transfinite algebra. *Bulletin of the American  
23 Mathematical Society*, 41(10), 667–670.