

# Dynamic Programming Deconstructed: Transformations of the Bellman Equation and Computational Efficiency<sup>1</sup>

Qingyin Ma<sup>a</sup> and John Stachurski<sup>b</sup>

<sup>a</sup>ISEM, Capital University of Economics and Business

<sup>b</sup>Research School of Economics, Australian National University

December 5, 2019

**ABSTRACT.** Some approaches to solving challenging dynamic programming problems, such as Q-learning, begin by transforming the Bellman equation into an alternative functional equation, in order to open up a new line of attack. Our paper studies this idea systematically, with a focus on boosting computational efficiency. We provide a characterization of the set of valid transformations of the Bellman equation, where validity means that the transformed Bellman equation maintains the link to optimality held by the original Bellman equation. We then examine the solutions of the transformed Bellman equations and analyze correspondingly transformed versions of the algorithms used to solve for optimal policies. These investigations yield new approaches to a variety of discrete time dynamic programming problems, including those with features such as recursive preferences or desire for robustness. Increased computational efficiency is demonstrated via time complexity arguments and numerical experiments.

*JEL Classifications:* C61, E00

*Keywords:* Dynamic programming, optimality, computational efficiency

## 1. INTRODUCTION

Dynamic programming is central to the analysis of intertemporal planning problems in management, operations research, economics, finance and other related disciplines (see, e.g., Bertsekas (2012)). When combined with statistical learning, dynamic programming also drives a number of strikingly powerful algorithms in artificial intelligence

---

<sup>1</sup>We thank Fedor Iskhakov, Takashi Kamihigashi, Larry Liu and Daisuke Oyama for valuable feedback and suggestions, as well as audience members at the Econometric Society meeting in Auckland in 2018 and the 2nd Conference on Structural Dynamic Models in Copenhagen in 2018. Financial support from ARC Discovery Grant DP120100321 is gratefully acknowledged.

*Email addresses:* qingyin.ma@cueb.edu.cn, john.stachurski@anu.edu.au

and automated decision systems (Kochenderfer (2015)). At the heart of dynamic programming lies the Bellman equation, a functional equation that summarizes the trade off between current and future rewards for a particular dynamic program. Under standard regularity conditions, solving the Bellman equation allows the researcher to assign optimal values to states and then compute optimal policies via Bellman’s principle of optimality (Bellman (1957)).

Despite its many successes, practical implementation of dynamic programming algorithms is often challenging, due to high dimensionality, irregularities and/or lack of data (see, e.g., Rust (1996)). A popular way to mitigate these problems is to try to change the angle of attack by rearranging the Bellman equation into an alternative functional form. One example of this is Q-learning, where the first step is to transform the Bellman equation into a new functional equation, the solution of which is the so-called Q-factor (see, e.g., Kochenderfer (2015) or Bertsekas (2012), Section 6.6.1). Working with the Q-factor turns out to be convenient when data on transition probabilities is scarce.

Other examples of transformations of the Bellman equation can be found in Kristensen et al. (2018), who investigate three versions of the Bellman equation associated with a fixed dynamic programming problem, corresponding to the value function, “expected value function” and “integrated value function” respectively. Similarly, Bertsekas (2012) discusses modifying the Bellman equation by integrating out “uncontrollable states” (Section 6.1.5). Each of these transformations of the Bellman equation creates new methods for solving for the optimal policy, since the transformations applied to the Bellman equation can be likewise applied to the iterative techniques used to solve the Bellman equation (e.g., value function iteration or policy iteration).

The purpose of this paper is twofold. First, we provide the first systematic investigation of these transformations, by developing a framework sufficiently general that all transformations of which we are aware can be expressed as special cases. We then examine the solutions of the transformed Bellman equations, along with correspondingly transformed versions of solution algorithms such as value function iteration or policy iteration. We provide a condition under which the transformed Bellman equations satisfy a version of Bellman’s principle of optimality, and the correspondingly transformed versions of the algorithms lead to optimal policies. This puts existing transformations on a firm theoretical foundation, in terms of their link to optimal policies, thereby eliminating the need to check optimality or convergence of algorithms on a case-by-case basis.

Second, on a more practical level, we use the framework developed above to create new transformations, or to apply existing transformations in new contexts, and to examine the convergence properties of the associated algorithms. Although there are many motivations for applying transformations to the Bellman equation (such as to facilitate learning, as in the case of Q-learning, or to simplify estimation as in [Rust \(1987\)](#)), the motivation that we focus on here is computational efficiency. In particular, we examine transformations of the Bellman equation that retain the links to optimality discussed above while reducing the dimension of the state space. Because the cost of high dimensionality is exponential, even small reductions in the number of variables can deliver speed gains of one or two orders of magnitude. These claims are verified theoretically and through numerical experiments when we consider applications.

One of the reasons that we are able to create new transformations is that we embed our theoretical results in a very general abstract dynamic programming framework, as found for example in [Bertsekas \(2013\)](#). This allows us to examine transformations of the Bellman equation in settings beyond the traditional additively separable case, which have been developed to better isolate risk preferences or accommodate notions of ambiguity. Examples of such objective functions can be found in [Iyengar \(2005\)](#), [Hansen and Sargent \(2008\)](#), [Ruszczyński \(2010\)](#) and [Bäuerle and Jaśkiewicz \(2018\)](#).

When analyzing solution methods and computational efficiency, we show that successive iterates of the transformed Bellman operator converge at the same rate as the original Bellman operator in the following sense: the  $n$ -th iterate of the Bellman operator can alternatively be produced by iterating the same number of times with the transformed Bellman operator and vice versa. This means that, to judge the relative efficiency, one only needs to consider the computational cost of a single iteration of the transformed Bellman operator versus the original Bellman operator.

When treating algorithms for computing optimal policies, we focus in particular on so-called optimistic policy iteration, which contains policy iteration and value function iteration as special cases, and which can typically be tuned to converge faster than either one in specific applications. We show that, under a combination of stability and monotonicity conditions, the sequence of policies generated by a transformed version of optimistic policy iteration, associated with a particular transformation of the Bellman equation, converges to optimality.

Some results related to ours and not previously mentioned can be found in [Rust \(1994\)](#), which discusses the connection between the fixed point of a transformed Bellman equation and optimality of the policy that results from choosing the maximal action at each state evaluated according to this fixed point. This is again a special case of what we

cover, specific to one specialized class of dynamic programming problems, with discrete choices and additively separable preferences, and refers to one specific transformation associated with the expected value function.

There is, of course, a very large literature on maximizing computational efficiency in solution methods for dynamic programming problems (see, e.g., [Powell \(2007\)](#)). The results presented here are, in many ways, complementary to this existing literature. For example, fitted value iteration is a popular technique for solving the Bellman equation via successive approximations, combined with a function approximation step (see, e.g., [Munos and Szepesvári \(2008\)](#)). This methodology could also, in principle, be applied to transformed versions of the Bellman equation and the successive approximation techniques for solving them examined in this paper.

The rest of our paper is structured as follows. Section 2 formulates the problem. Section 3 discusses optimality and algorithms. Section 4 gives a range of applications. Longer proofs are deferred to the appendix.

## 2. GENERAL FORMULATION

This section presents an abstract dynamic programming problem and the key concepts and operators related to the formulation.

**2.1. Problem Statement.** We use  $\mathbb{N}$  to denote the set of natural numbers and  $\mathbb{N}_0 := \{0\} \cup \mathbb{N}$ , while  $\mathbb{R}^E$  is the set of real-valued functions defined on some set  $E$ . In what follows, a dynamic programming problem consists of

- a nonempty set  $\mathbf{X}$  called the *state space*,
- a nonempty set  $\mathbf{A}$  called the *action space*,
- a nonempty correspondence  $\Gamma$  from  $\mathbf{X}$  to  $\mathbf{A}$  called the *feasible correspondence*, along with the associated set of *state-action pairs*

$$\mathbf{F} := \{(x, a) \in \mathbf{X} \times \mathbf{A} : a \in \Gamma(x)\},$$

- a subset  $\mathcal{V}$  of  $\mathbb{R}^{\mathbf{X}}$  called the set of *candidate value functions*, and
- a *state-action aggregator*  $H$  mapping  $\mathbf{F} \times \mathcal{V}$  to  $\mathbb{R} \cup \{-\infty\}$ .

The interpretation of  $H(x, a, v)$  is the lifetime value associated with choosing action  $a$  at current state  $x$  and then continuing with a reward function  $v$  attributing value to states. In other words,  $H(x, a, v)$  is an abstract representation of the value to be maximized on the right hand side of the Bellman equation (see (1) below). This abstract representation accommodates both additively separable and nonseparable preferences

and is based on Bertsekas (2013). The sets  $\mathbf{X}$  and  $\mathbf{A}$  are, at this point, arbitrary. They can, for example, be finite or subsets of Euclidean vector space.

We associate to this abstract dynamic program the Bellman equation

$$v(x) = \sup_{a \in \Gamma(x)} H(x, a, v) \quad \text{for all } x \in \mathbf{X}. \quad (1)$$

Stating that  $v \in \mathcal{V}$  solves the Bellman equation is equivalent to stating that  $v$  is a fixed point of the Bellman operator, which we denote by  $T$  and define by

$$Tv(x) = \sup_{a \in \Gamma(x)} H(x, a, v) \quad (x \in \mathbf{X}, v \in \mathcal{V}). \quad (2)$$

**Example 2.1.** In a traditional infinite horizon finite state Markov decision process, an agent observes a state  $x$  in finite set  $\mathbf{X}$  and responds with action  $a$  from  $\Gamma(x)$ , a subset of finite set  $\mathbf{A}$ . The state then updates to  $x'$  next period with probability  $p(x, a, x')$ . The objective is to choose a policy  $\sigma: \mathbf{X} \rightarrow \mathbf{A}$  such that, when  $\{a_t\}$  is selected according to  $a_t = \sigma(x_t)$  at each  $t$ , the objective  $\mathbb{E} \sum_{t \geq 0} \beta^t r(x_t, a_t)$  is maximized. Here  $\beta \in (0, 1)$  is a discount factor and  $r$  is a real-valued reward function. This fits directly into our framework if we set  $\mathcal{V}$  to be all functions from  $\mathbf{X}$  to  $\mathbb{R}$  and

$$H(x, a, v) = r(x, a) + \beta \sum_{x' \in \mathbf{X}} v(x') p(x, a, x'). \quad (3)$$

In particular, inserting the right hand side of (3) into (1) reproduces the standard Bellman equation for this problem (see, e.g., Bertsekas (2013), Example 1.2.2).

More detailed examples, involving those with nonseparable preferences, are given in Section 2.4 and Section 4 below.

**2.2. Policies and Assumptions.** Let  $\Sigma$  denote the set of *feasible policies*, which we define as all  $\sigma: \mathbf{X} \rightarrow \mathbf{A}$  satisfying  $\sigma(x) \in \Gamma(x)$  for all  $x \in \mathbf{X}$  and the regularity condition

$$v \in \mathcal{V} \text{ and } w(x) = H(x, \sigma(x), v) \text{ on } \mathbf{X} \implies w \in \mathcal{V}. \quad (4)$$

Given  $v \in \mathcal{V}$ , a feasible policy  $\sigma$  with the property that

$$H(x, \sigma(x), v) = \sup_{a \in \Gamma(x)} H(x, a, v) \quad \text{for all } x \in \mathbf{X} \quad (5)$$

is called *v-greedy*. In other words, a *v-greedy* policy is one we obtain by treating  $v$  as the value function and maximizing the right hand side of the Bellman equation.

**Assumption 2.1.** There exists a subset  $\mathbb{V}$  of  $\mathcal{V}$  such that  $T$  maps elements of  $\mathbb{V}$  into itself, and a *v-greedy* policy exists in  $\Sigma$  for each  $v$  in  $\mathbb{V}$ .

Assumption 2.1 guarantees the existence of stationary policies and allows us to work with maximal decision rules instead of suprema. It is obviously satisfied for any dynamic program where the set of actions is finite, such as Example 2.1 (by letting  $\mathbb{V} = \mathcal{V}$ ), and extends to infinite choice problems when primitives are sufficiently continuous and the choice set at each action is compact.

Given  $\sigma \in \Sigma$ , any function  $v_\sigma$  in  $\mathcal{V}$  satisfying

$$v_\sigma(x) = H(x, \sigma(x), v_\sigma) \quad \text{for all } x \in \mathsf{X}$$

is called a  $\sigma$ -value function. We understand  $v_\sigma(x)$  as the lifetime value of following policy  $\sigma$  now and forever, starting from current state  $x$ .

**Assumption 2.2.** For each  $\sigma \in \Sigma$ , there is exactly one  $\sigma$ -value function in  $\mathcal{V}$ , denoted by  $v_\sigma$ .

Assumption 2.2 is required for our optimization problem to be well defined, and is easy to verify for regular Markov decision problems such as the one described in Example 2.1 (see, e.g., Bertsekas (2013), Example 1.2.1). When rewards are unbounded or dynamic preferences are specified recursively, the restrictions on primitives used to establish this assumption vary substantially across applications (see, e.g., Marinacci and Montrucchio (2010), Bäuerle and Jaśkiewicz (2018), or Bloise and Vailakis (2018) for sufficient conditions in these contexts, and Bertsekas (2013) for general discussion).

**2.3. Decompositions and Plan Factorizations.** Next we introduce plan factorizations (which correspond to transformations of the Bellman equation) in a relatively abstract form, motivated by the desire to accommodate existing transformations observed in the literature and admit new ones. As a start, let  $\mathcal{G}$  (resp.,  $\mathbb{G}$ ) be the set of functions  $g$  in  $\mathbb{R}^{\mathsf{F}}$  such that  $g = W_0 v$  for some  $v \in \mathcal{V}$  (resp., for some  $v \in \mathbb{V}$ ). Let  $\mathcal{H}$  (resp.,  $\mathbb{H}$ ) be all functions  $h$  in  $\mathbb{R}^{\mathsf{F}}$  such that, for some  $v \in \mathcal{V}$  (resp., for some  $v \in \mathbb{V}$ ), we have  $h(x, a) = H(x, a, v)$  for all  $(x, a) \in \mathsf{F}$ . Let  $M$  be the operator defined at  $h \in \mathcal{H}$  by

$$(Mh)(x) = \sup_{a \in \Gamma(x)} h(x, a). \tag{6}$$

The operator  $M$  maps elements of  $\mathbb{H}$  into  $\mathbb{V}$ . Indeed, given  $h \in \mathbb{H}$ , there exists by definition a  $v \in \mathbb{V}$  such that  $h(x, a) = H(x, a, v)$  at each  $(x, a) \in \mathsf{F}$ . We then have  $Mh = Tv \in \mathbb{V}$  by (2) and Assumption 2.1.

A *plan factorization* associated with the dynamic program described above is a pair of operators  $(W_0, W_1)$  such that

- (1)  $W_0$  is defined on  $\mathcal{V}$  and takes values in  $\mathcal{G}$ ,

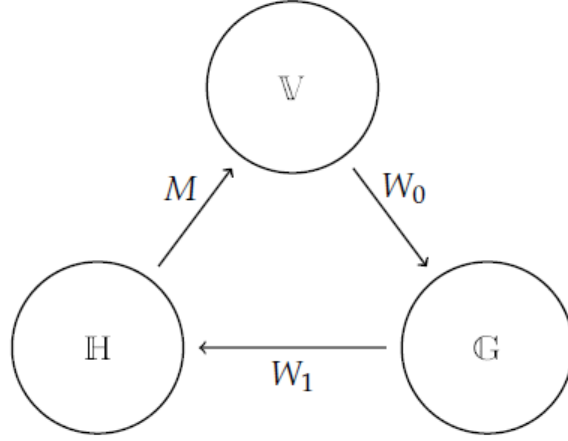


FIGURE 1. The one-shift operators

- (2)  $W_1$  is defined on  $\mathcal{G}$  and takes values in  $\mathcal{H}$ , and
- (3) the operator composition  $W_1 W_0$  satisfies

$$(W_1 W_0 v)(x, a) = H(x, a, v) \quad \text{for all } (x, a) \in \mathcal{F}, v \in \mathcal{V}. \quad (7)$$

Equation (7) states that  $W_0$  and  $W_1$  provide a decomposition (or factorization) of  $H$ , so that each element  $W_i$  implements a separate component of the two stage (i.e., present and future) planning problem associated with the Bellman equation.

Given the definition of  $M$  in (6) and the factorization requirement (7), the Bellman operator  $T$  from (2) can now be expressed as

$$T = M W_1 W_0. \quad (8)$$

In the following, when we discuss the Bellman operator (and the refactored Bellman operator to be defined below), we confine the domain of the operators  $W_0$ ,  $W_1$  and  $M$  to  $\mathbb{V}$ ,  $\mathbb{G}$  and  $\mathbb{H}$ , respectively. In that circumstance,  $T$  is a cycle starting from  $\mathbb{V}$ . In particular,  $W_0$  maps  $\mathbb{V}$  into  $\mathbb{G}$ , and  $W_1$  maps  $\mathbb{G}$  into  $\mathbb{H}$  by definition, while  $M$  maps  $\mathbb{H}$  into  $\mathbb{V}$  as has been shown above. A visualization is given in Figure 1.

Corresponding to this same plan factorization  $(W_0, W_1)$ , we introduce a *refactored Bellman operator*  $S$  on  $\mathbb{G}$  defined by

$$S = W_0 M W_1. \quad (9)$$

In terms of Figure 1,  $S$  is a cycle starting from  $\mathbb{G}$ . Corresponding to  $S$ , we have the *refactored Bellman equation*

$$g(x, a) = (W_0 M W_1 g)(x, a) \quad \text{for all } (x, a) \in \mathbb{F}. \quad (10)$$

The transformations of the Bellman equation we wish to consider all correspond to a version of (10), under some suitable specification of  $W_0$  and  $W_1$ .

**2.4. Examples.** In this section we briefly illustrate the framework with examples. The examples contain both common transformations and new ones. Throughout,  $\beta \in (0, 1)$  denotes the discount factor.

**2.4.1.  $Q$ -Factors.** Consider the traditional finite state Markov decision process from Example 2.1, where  $\mathcal{V}$  is all functions from  $\mathbf{X}$  to  $\mathbb{R}$  and  $H$  is as given in (3). Let  $\mathbb{V} = \mathcal{V}$ . One special case of the refactored Bellman equation is the equation for determining optimal  $Q$ -factors that forms the basis of  $Q$ -learning. If  $I$  is the identity map and we set

$$(W_0 v)(x, a) = r(x, a) + \beta \sum_{x' \in \mathbf{X}} v(x') p(x, a, x') \quad \text{and} \quad W_1 = I,$$

then  $(W_0, W_1)$  is a plan factorization, since  $(W_1 W_0 v)(x, a) = H(x, a, v)$  when the latter is given by (3). For this plan factorization, the refactored Bellman operator  $S$  from (9) becomes  $S = W_0 M$ , or, more explicitly,

$$(Sg)(x, a) = r(x, a) + \beta \sum_{x' \in \mathbf{X}} \max_{a' \in \Gamma(x')} g(x', a') p(x, a, x').$$

Solving for  $g = Sg$  is exactly the problem associated with computation of optimal  $Q$ -factors. See, for example, Equation (1.4) of Bertsekas and Yu (2012) or Equation (2.24) of Bertsekas (2012). Thus, the underlying functional equations and fixed point problems of  $Q$ -learning correspond to the special case of our theory where  $W_1$  is the identity.

**2.4.2. The Standard Case.** The previous example showed how the problem of obtaining  $Q$ -factors is one kind of plan factorization—in fact an extreme kind, where the element  $W_1$  of the pair  $(W_0, W_1)$  is the identity. The other extreme case is when  $W_0$  is the identity. Then, the standard and refactored operators  $T$  and  $S$  coincide, as can be confirmed by consulting (8) and (9). Thus, the standard Bellman equation and the  $Q$ -factor equation can be viewed as the two extremal cases of plan factorization.



2.4.3. *Expected Value Transformation.* Continuing with the Markov decision process from Example 2.1, where  $H$  is as given in (3), consider the plan factorization  $(W_0, W_1)$  defined by

$$(W_0 v)(x, a) = \sum_{x' \in \mathbf{X}} v(x') p(x, a, x') \quad \text{and} \quad (W_1 g)(x, a) = r(x, a) + \beta g(x, a). \quad (11)$$

Evidently (11) yields a plan factorization, in the sense that  $(W_1 W_0 v)(x, a) = H(x, a, v)$  when the latter is given by (3). For this plan factorization, the refactored Bellman operator  $S$  from (9) becomes

$$(Sg)(x, a) = (W_0 M W_1 g)(x, a) = \sum_{x' \in \mathbf{X}} \max_{a' \in \Gamma(x')} \{r(x', a') + \beta g(x', a')\} p(x, a, x'). \quad (12)$$

With this choice of plan factorization, the operator  $S$  is a generalization of the “expected value operator” used by Rust (1987) in the context of optimal timing of decisions for bus engine replacement. At this point in time it is not apparent why one should favor the use of  $S$  over the regular Bellman operator  $T$ , and indeed there might be no advantage. The benefit of having a general optimality theory based around  $S$  rather than  $T$ , as constructed below, is the option of using  $S$  when the structure of the problem implies that doing so *is* advantageous. Examples are shown in Section 4.

2.4.4. *Optimal Stopping.* Consider an optimal stopping problem with additively separable preferences (see, e.g., Monahan (1980) or Peskir and Shiryaev (2006)), where an agent observes current state  $x$  and chooses between stopping and continuing. Stopping generates terminal reward  $r(x)$  while continuing yields flow continuation reward  $c(x)$ . If the agent continues, the next period state  $x'$  is drawn from  $P(x, dx')$  and the process repeats. The agent takes actions in order to maximize expected lifetime rewards. The state-action aggregator  $H$  for this model is

$$H(x, a, v) = a r(x) + (1 - a) \left[ c(x) + \beta \int v(x') P(x, dx') \right], \quad (13)$$

where  $a \in \Gamma(x) := \{0, 1\}$  is a binary choice variable with  $a = 1$  indicating the decision to stop. Let  $\mathbf{X}$  be a Borel subset of  $\mathbb{R}^m$ ,  $r$  and  $c$  be bounded and continuous,  $\mathcal{V} = \mathbb{V}$  be the set of bounded continuous functions on  $\mathbf{X}$ , and  $P$  satisfy the Feller property (i.e.,  $x \mapsto \int v(x') P(x, dx')$  is bounded and continuous whenever  $v$  is). Assumption 2.1 holds because the action space is discrete and Assumption 2.2 holds by a standard contraction argument (see, e.g., Section 2.1 of Bertsekas (2013)).

One possible plan factorization is

$$(W_0 v)(x) = \int v(x') P(x, dx') \quad \text{and} \quad (W_1 g)(x) = a r(x) + (1 - a) [c(x) + \beta g(x)].$$

The refactored Bellman operator  $S = W_0 M W_1$  is then given by

$$(Sg)(x) = \int \max \{r(x'), c(x') + \beta g(x')\} P(x, dx'). \quad (14)$$

There exist applications in the field of optimal stopping where iteration with the operator  $S$  given in (14) is computationally more efficient than iterating with the standard Bellman operator  $T$  (see, e.g., Rust (1987)). Section 4.2 expands on this point.

2.4.5. *Risk Sensitive Preferences.* Consider a risk-sensitive control problem where  $H$  can be expressed as

$$H(x, a, v) = \left\{ r(x, a) - \frac{\beta}{\gamma} \log \int \exp[-\gamma v(x')] P(x, a, dx') \right\}. \quad (15)$$

Here  $P(x, a, dx')$  gives one-step state transition probabilities given state  $x \in \mathbb{R}_+$  and action  $a \in [0, x]$ , while  $\gamma > 0$  controls the degree of risk sensitivity. One example of this set up occurs in Bäuerle and Jaśkiewicz (2018), where  $H(x, a, v)$  has the interpretation of lifetime rewards when an agent owns current capital  $x$ , choses current consumption  $a$ , receives current reward  $r(x, a)$  and then transitions to the next period with capital drawn from  $P(x, a, dx')$ . Sufficient conditions for Assumptions 2.1–2.2 can be found in Theorem 1 of that paper.

There are many plan factorizations we could consider here, one of which extends the Q-factor transformation in Section 2.4.1, and another of which extends the expected value transformation in Section 2.4.3. Taking the latter as an example, we set

$$(W_0 v)(x, a) = -\frac{\beta}{\gamma} \log \int \exp[-\gamma v(x')] P(x, a, dx') \\ \text{and } (W_1 g)(x, a) = r(x, a) + g(x, a). \quad (16)$$

Evidently (16) yields a plan factorization, in the sense that  $(W_1 W_0 v)(x, a) = H(x, a, v)$  when the latter is given by (15). The refactored Bellman operator  $S = W_0 M W_1$  is given by

$$(Sg)(x, a) = -\frac{\beta}{\gamma} \log \int \exp \left[ -\gamma \max_{a' \in \Gamma(x')} \{r(x', a') + g(x', a')\} \right] P(x, a, dx'). \quad (17)$$

2.5. **Refactored Policy Values.** Returning to the general setting, we also wish to consider the value of individual policies under the transformation associated with a given plan factorization  $(W_0, W_1)$ . This will be important when we consider topics such as policy function iteration. To this end, for each  $\sigma \in \Sigma$  and each  $h \in \mathcal{H}$ , we define the operator  $M_\sigma: \mathcal{H} \rightarrow \mathcal{V}$  by

$$M_\sigma h(x) := h(x, \sigma(x)). \quad (18)$$

That  $M_\sigma$  does in fact map  $\mathcal{H}$  into  $\mathcal{V}$  follows from the definition of these spaces and (4). In particular, if  $h \in \mathcal{H}$ , by definition, there exists a  $v \in \mathcal{V}$  such that  $h(x, a) = H(x, a, v)$  for all  $(x, a) \in \mathbf{F}$ . Then  $M_\sigma h \in \mathcal{V}$  follows directly from (4). Comparing with  $M$  defined in (6), we have

$$M_\sigma h \leq Mh \quad \text{for each } h \in \mathcal{H} \text{ and } \sigma \in \Sigma. \quad (19)$$

Given  $\sigma \in \Sigma$ , the operator  $T_\sigma$  from  $\mathcal{V}$  to itself defined by  $T_\sigma v(x) = H(x, \sigma(x), v)$  for all  $x \in \mathbf{X}$  will be called the  $\sigma$ -value operator. By construction, it has the property that  $v_\sigma$  is the  $\sigma$ -value function corresponding to  $\sigma$  if and only if it is a fixed point of  $T_\sigma$ . With the notation introduced above, we can express it as

$$T_\sigma = M_\sigma W_1 W_0. \quad (20)$$

Note that  $T_\sigma$  is a cycle starting from  $\mathcal{V}$ . In particular,  $W_0$  maps  $\mathcal{V}$  into  $\mathcal{G}$ , and  $W_1$  maps  $\mathcal{G}$  into  $\mathcal{H}$  by the definition of plan factorization, while  $M_\sigma$  maps  $\mathcal{H}$  into  $\mathcal{V}$  as shown above.

Analogous to the definition of the refactored Bellman operator in (9), we introduce the *refactored  $\sigma$ -value operator* on  $\mathcal{G}$

$$S_\sigma := W_0 M_\sigma W_1 \quad (21)$$

corresponding to a given plan factorization  $(W_0, W_1)$ . Similarly,  $S_\sigma$  is a cycle starting from  $\mathcal{G}$ . A fixed point  $g_\sigma$  of  $S_\sigma$  is called a *refactored  $\sigma$ -value function*. The value  $g_\sigma(x, a)$  can be interpreted as the value of following policy  $\sigma$  in all subsequent periods, conditional on current action  $a$ .

**Example 2.2.** Recall the expected value function factorization of a finite state Markov decision process discussed in Section 2.4.3. For a given policy  $\sigma$ , the refactored  $\sigma$ -value operator can be obtained by replacing  $M$  in (12) with  $M_\sigma$ , which yields

$$\begin{aligned} (S_\sigma g)(x, a) &= (W_0 M_\sigma W_1 g)(x, a) \\ &= \sum_{x' \in \mathbf{X}} \{r(x', \sigma(x')) + \beta g(x', \sigma(x'))\} p(x, a, x'). \end{aligned} \quad (22)$$

**Example 2.3.** Recall the plan factorization (14) applied to the optimal stopping problem with aggregator  $H$  given by (13). A feasible policy is a map  $\sigma$  from  $\mathbf{X}$  to  $\mathbf{A} = \{0, 1\}$ . The refactored  $\sigma$ -value operator corresponding to this plan factorization is

$$(S_\sigma g)(x) = \int \{\sigma(x')r(x') + (1 - \sigma(x'))[c(x') + \beta g(x')]\} P(x, dx'). \quad (23)$$

**2.6. Fixed Points and Iteration.** We begin with some preliminary results on properties of the operators defined above. Our first result states that, to iterate with  $T$ , one can alternatively iterate with  $S$ , since iterates of  $S$  can be converted directly into iterates of  $T$ . Moreover, the converse is also true, and similar statements hold for the policy operators:

**Proposition 2.1.** *For every  $n \in \mathbb{N}$ , we have*

$$S^n = W_0 T^{n-1} M W_1 \quad \text{and} \quad T^n = M W_1 S^{n-1} W_0.$$

*Moreover, for every  $n \in \mathbb{N}$  and every  $\sigma \in \Sigma$ , we have*

$$S_\sigma^n = W_0 T_\sigma^{n-1} M_\sigma W_1 \quad \text{and} \quad T_\sigma^n = M_\sigma W_1 S_\sigma^{n-1} W_0.$$

Proposition 2.1 follows from the definitions of  $T$  and  $S$  along with a simple induction argument. While the proof is relatively straightforward, the result is important because that it provides a metric-free statement of the fact that the sequence of iterates of any refactored Bellman operator converges at the same rate as that of the standard Bellman operator.

The next result shows a fundamental connection between the two forms of the Bellman operator in terms of their fixed points.

**Proposition 2.2.** *The Bellman operator  $T$  admits a unique fixed point  $\bar{v}$  in  $\mathbb{V}$  if and only if the refactored Bellman operator  $S$  admits a unique fixed point  $\bar{g}$  in  $\mathbb{G}$ . Whenever these fixed points exist, they are related by  $\bar{v} = M W_1 \bar{g}$  and  $\bar{g} = W_0 \bar{v}$ .*

Thus, if a unique fixed point of the Bellman operator is desired but  $T$  is difficult to work with, a viable option is to show that  $S$  has a unique fixed point, compute it, and then recover the unique fixed point of  $T$  via  $\bar{v} = M W_1 \bar{g}$ .

A result analogous to Proposition 2.2 holds for the policy operators:

**Proposition 2.3.** *Given  $\sigma \in \Sigma$ , the refactored  $\sigma$ -value operator  $S_\sigma$  has a unique fixed point  $g_\sigma$  in  $\mathcal{G}$  if and only if  $T_\sigma$  has a unique fixed point  $v_\sigma$  in  $\mathcal{V}$ . Whenever these fixed points exist, they are related by  $v_\sigma = M_\sigma W_1 g_\sigma$  and  $g_\sigma = W_0 v_\sigma$ .*

Proposition 2.3 implies that, under Assumption 2.2, there exists exactly one refactored  $\sigma$ -value function  $g_\sigma$  in  $\mathcal{G}$  for every  $\sigma \in \Sigma$ . Proposition 2.3 is also useful for the converse implication. In particular, to establish Assumption 2.2, which is often nontrivial when preferences are not additively separable, we can equivalently show that  $S_\sigma$  has a unique fixed point in  $\mathcal{G}$ .

### 3. OPTIMALITY

Next we turn to optimality, with a particular focus on the properties that must be placed on a given plan factorization in order for the corresponding (refactored) Bellman equation to lead us to optimal actions. Our first step, however, is to define optimality and recall some standard results.

**3.1. Fixed Points and Optimal Policies.** The *value function* associated with our dynamic program is defined at  $x \in \mathbf{X}$  by

$$v^*(x) = \sup_{\sigma \in \Sigma} v_\sigma(x). \quad (24)$$

A feasible policy  $\sigma^*$  is called *optimal* if  $v_{\sigma^*} = v^*$  on  $\mathbf{X}$ . Given  $\sigma \in \Sigma$ , *Bellman’s principle of optimality* states that

$$\sigma \text{ is an optimal policy} \iff \sigma \text{ is } v^*\text{-greedy.} \quad (25)$$

The next theorem is a simple extension of foundational optimality results for dynamic decision problems that link the Bellman equation to optimality (see, e.g., Bertsekas (2013)). Its proof is omitted. The assumptions of Section 2.2 are taken to be in force.

**Theorem 3.1.** *The next two statements are equivalent:*

- (1) *The value function  $v^*$  lies in  $\mathbb{V}$  and satisfies the Bellman equation (1).*
- (2) *Bellman’s principle of optimality (25) holds and the set of optimal policies is nonempty.*

The motivation behind the transformations we consider in this paper is that the Bellman equation can be refactored into a more convenient form without affecting optimality. Thus, it is natural to ask when—and under what conditions—a result analogous to Theorem 3.1 holds for the refactored Bellman equation (10). We will show that an analogous result obtains whenever a form of monotonicity holds for the transformation being used in the plan factorization.

Before we get to this result, however, we need to address the following complication: there are two distinct functions that can take the part of  $v^*$  in Theorem 3.1. One is the “rotated value function”

$$\hat{g} := W_0 v^*. \quad (26)$$

The second is

$$g^*(x, a) := \sup_{\sigma \in \Sigma} g_\sigma(x, a). \quad (27)$$

We call  $g^*$  the *refactored value function*. The definition of  $g^*$  directly parallels the definition of the value function in (24), with  $g^*(x, a)$  representing the maximal value that can be obtained from state  $x$  conditional on choosing  $a$  in the current period. (Since Assumption 2.2 is in force, the set of functions  $\{g_\sigma\}$  in the definition of  $g^*$  is well defined by Proposition 2.3, and hence so is  $g^*$  as an extended real-valued function, although it might or might not live in  $\mathcal{G}$ .)

As shown below, the functions  $\hat{g}$  and  $g^*$  are not in general equal, although they become so when a certain form of monotonicity is imposed. Moreover, under this same monotonicity condition, if one and hence both of these functions are fixed points of  $S$ , we obtain valuable optimality results.

In stating these results, we recall that a map  $A$  from one partially ordered set  $(E, \leq)$  to another such set  $(F, \leq)$  is called *monotone* if  $x \leq y$  implies  $Ax \leq Ay$ . Below, monotonicity is with respect to the pointwise partial order on  $\mathbb{V}$ ,  $\mathbb{G}$  and  $\mathbb{H}$ . Moreover, given  $g \in \mathbb{G}$ , a policy  $\sigma \in \Sigma$  is called  *$g$ -greedy* if  $M_\sigma W_1 g = M W_1 g$ . At least one  $g$ -greedy policy exists for every  $g \in \mathbb{G}$ . Indeed,  $g \in \mathbb{G}$  implies the existence of a  $v \in \mathbb{V}$  such that  $g = W_0 v$ , and to this  $v$  there corresponds a  $v$ -greedy policy  $\sigma$  by Assumption 2.1. At this  $\sigma$  we have  $H(x, \sigma(x), v) = \sup_{a \in \Gamma(x)} H(x, a, v)$  for every  $x \in \mathbf{X}$ , or, equivalently,  $M_\sigma W_1 W_0 v = M W_1 W_0 v$  pointwise on  $\mathbf{X}$ . Since  $g = W_0 v$ , this policy is  $g$ -greedy.

We can now state our first significant optimality result. It shows that, under some monotonicity conditions, the standard Bellman equation and the refactored Bellman equation have “equal rights,” in the sense that both can be used to obtain the same optimal policy, and both satisfy a version of Bellman’s principle of optimality.

**Theorem 3.2.** *Let  $(W_0, W_1)$  be a plan factorization, let  $\hat{g}$  be as defined in (26) and let  $g^*$  be as defined in (27). If  $W_0$  and  $W_1$  are both monotone, then the following statements are equivalent:*

- (1)  $g^*$  lies in  $\mathbb{G}$  and satisfies the refactored Bellman equation (10).
- (2)  $v^*$  lies in  $\mathbb{V}$  and satisfies the Bellman equation (1).

*If these conditions hold, then  $g^* = \hat{g}$ , the set of optimal policies is nonempty and, for  $\sigma \in \Sigma$ , we have*

$$\sigma \text{ is an optimal policy} \iff \sigma \text{ is } g^*\text{-greedy} \iff \sigma \text{ is } v^*\text{-greedy.} \quad (28)$$

The monotonicity conditions on  $W_0$  and  $W_1$  clearly hold in all of the example transformations discussed in Sections 2.4.1–2.4.4. It is possible to envisage settings where

they fail, however. For example, if we replace the factorization  $(W_0, W_0)$  in (16) with

$$(W_0 v)(x, a) = \int \exp[-\gamma v(x')] P(x, a, dx')$$

$$\text{and } (W_1 g)(x, a) = r(x, a) - \frac{\beta}{\gamma} \log g(x, a),$$

we again have a valid plan factorization (since  $(W_1 W_0 v)(x, a) = H(x, a, v)$  when the latter is given by (15)) but neither  $W_0$  nor  $W_1$  is monotone.

In general, monotonicity of  $W_0$  and  $W_1$  cannot be dropped from Theorem 3.2 without changing its conclusions. In Section 6.4 of the Appendix we exhibit dynamic programs and plan factorizations that illustrate the following possibilities:

- (1)  $\hat{g} \neq g^*$ .
- (2)  $Tv^* = v^*$  and yet  $Sg^* \neq g^*$ .
- (3)  $Sg^* = g^*$  and yet  $Tv^* \neq v^*$ .

**3.2. Sufficient Conditions.** Theorem 3.2 tells us that if the stated monotonicity condition holds and  $Sg^* = g^*$ , then we can be assured of the existence of optimal policies and have a means to characterize them. What we lack is a set of sufficient conditions under which the refactored Bellman operator has a unique fixed point that is equal to  $g^*$ . The theorem stated in this section fills that gap.

To state the theorem, we recall that a self-mapping  $A$  on a topological space  $U$  is called *asymptotically stable* on  $U$  if  $A$  has a unique fixed point  $\bar{u}$  in  $U$  and  $A^n u \rightarrow \bar{u}$  as  $n \rightarrow \infty$  for all  $u \in U$ . In the following, we assume that there exists a Hausdorff topology  $\tau$  on  $\mathcal{G}$  under which the pointwise partial order is closed (i.e., its graph is closed in the product topology on  $\mathcal{G} \times \mathcal{G}$ ). The key implication is that if  $f_n \rightarrow f$  and  $g_n \rightarrow g$  under  $\tau$  and  $f_n \leq g_n$  for all  $n$ , then  $f \leq g$ . Asymptotic stability is with respect to this topology.

In the theorem below,  $(W_0, W_1)$  is a given plan factorization and  $g^*$  is the refactored value function.

**Theorem 3.3.** *If  $W_0$  and  $W_1$  are monotone,  $S$  is asymptotically stable on  $\mathbb{G}$  and  $\{S_\sigma\}_{\sigma \in \Sigma}$  are asymptotically stable on  $\mathcal{G}$ , then*

- (1)  $g^*$  is the unique solution to the refactored Bellman equation in  $\mathbb{G}$ ,
- (2)  $\lim_{k \rightarrow \infty} S^k g = g^*$  under  $\tau$  for all  $g \in \mathbb{G}$ ,
- (3) at least one optimal policy exists, and
- (4) a feasible policy is optimal if and only if it is  $g^*$ -greedy.

In Theorem 3.3 we eschewed a contractivity assumption on  $S$  or  $S_\sigma$ , since such an assumption is problematic in certain applications (see, e.g., Bertsekas (2013)). Nevertheless, contraction methods can be applied to many other problems. For this reason we add the next proposition, which is a refactored analog of the classical theory of dynamic programming based around contraction mappings.

Let  $\|\cdot\|_\kappa$  be defined at  $f \in \mathbb{R}^F$  by  $\|f\|_\kappa = \sup |f/\kappa|$ . Here the supremum is over all  $(x, a) \in F$  and  $\kappa$  is a fixed “weighting function”, i.e., an element of  $\mathbb{R}^F$  satisfying  $\inf \kappa > 0$  on  $F$ . This mapping defines a norm on  $b_\kappa F$ , the set of  $f \in \mathbb{R}^F$  such that  $\|f\|_\kappa < \infty$ .

**Proposition 3.4.** *Let  $W_0$  and  $W_1$  be monotone. If  $\mathcal{G}$  and  $\mathbb{G}$  are closed subsets of  $b_\kappa F$  and there exists a positive constant  $\alpha$  such that  $\alpha < 1$  and*

$$\|S_\sigma g - S_\sigma g'\|_\kappa \leq \alpha \|g - g'\|_\kappa \quad \text{for all } g, g' \in \mathcal{G} \text{ and } \sigma \in \Sigma, \quad (29)$$

*then  $S$  is a contraction mapping on  $(\mathbb{G}, \|\cdot\|_\kappa)$  of modulus  $\alpha$ ,  $S$  is asymptotically stable on  $\mathbb{G}$ , and the conclusions of Theorem 3.3 all hold.*

The purpose of  $\kappa$  here is to control for unbounded rewards (see, e.g., Bertsekas (2013)). When rewards are bounded we can take  $\kappa \equiv 1$ , in which case  $\|\cdot\|_\kappa$  is the ordinary supremum norm and  $b_\kappa F$  is just the bounded functions on  $F$ . Notice also that the contractivity requirement in (29) is imposed on  $S_\sigma$  rather than  $S$ , making the condition easier to verify (since  $S_\sigma$  does not involve a maximization step).

**3.3. Policy Iteration.** There are algorithms besides value function iteration that achieve faster convergence in some applications. These include Howard’s policy iteration algorithm and a popular variant called optimistic policy iteration (see, e.g., Puterman and Shin (1982) or Bertsekas (2012)). Optimistic policy iteration is important because it includes value function iteration and Howard’s policy iteration algorithm as limiting cases. In this section we give conditions under which optimistic policy iteration is successful in the setting of a given plan factorization  $(W_0, W_1)$ . We make the following assumption.

**Assumption 3.1.** At least one  $v$ -greedy policy exists in  $\Sigma$  for each  $v$  in  $\mathcal{V}$ .

Note that Assumption 3.1 implies Assumption 2.1. In particular, by Assumption 3.1, for each  $v \in \mathcal{V}$ , there exists a  $\sigma$  in  $\Sigma$  such that  $Tv(x) = H(x, \sigma(x), v)$  for all  $x \in X$ . (4) then implies that  $Tv \in \mathcal{V}$ . Therefore,  $T$  is a self-map on  $\mathcal{V}$  and Assumption 2.1 holds by letting  $\mathbb{V} = \mathcal{V}$ .



The standard algorithm starts with an initial candidate  $v_0 \in \mathcal{V}$  and generates sequences  $\{\sigma_k^v\}$ ,  $\{\Sigma_k^v\}$  in  $\Sigma$  and  $\{v_k\}$  in  $\mathcal{V}$  by taking

$$\sigma_k^v \in \Sigma_k^v \quad \text{and} \quad v_{k+1} = T_{\sigma_k^v}^{m_k} v_k \quad \text{for all } k \in \mathbb{N}_0, \quad (30)$$

where  $\Sigma_k^v$  is the set of  $v_k$ -greedy policies, and  $\{m_k\}$  is a sequence of positive integers. The first step of equation (30) is called *policy improvement*, while the second step is called *partial policy evaluation*. If  $m_k = 1$  for all  $k$  then the algorithm reduces to value function iteration.

To extend this idea to the refactored case, we take an initial candidate  $g_0 \in \mathcal{G}$  and generate sequences  $\{\sigma_k^g\}$ ,  $\{\Sigma_k^g\}$  in  $\Sigma$  and  $\{g_k\}$  in  $\mathcal{G}$  via

$$\sigma_k^g \in \Sigma_k^g \quad \text{and} \quad g_{k+1} = S_{\sigma_k^g}^{m_k} g_k \quad \text{for all } k \in \mathbb{N}_0, \quad (31)$$

where  $\Sigma_k^g$  is the set of  $g_k$ -greedy policies. The next result shows that (30) and (31) indeed generate the same sequences of greedy policies.

**Theorem 3.5.** *If  $v_0 \in \mathcal{V}$  and  $g_0 = W_0 v_0$ , then sequences  $\{\sigma_k\}$  and  $\{\Sigma_k\}$  in  $\Sigma$  satisfy (30) if and only if they satisfy (31). Moreover, the generated  $\{v_k\}$  and  $\{g_k\}$  sequences satisfy  $g_k = W_0 v_k$  for all  $k$ .*

Moreover, optimistic policy iteration via the refactored Bellman operator converges, as for the standard Bellman operator:

**Theorem 3.6.** *Let  $W_0$  and  $W_1$  be monotone and let  $S$  and  $\{S_\sigma\}_{\sigma \in \Sigma}$  be asymptotically stable on  $\mathcal{G}$ . If  $g_0 \leq Sg_0$ , then  $g_k \leq g_{k+1}$  for all  $k$  and  $g_k \rightarrow g^*$  as  $k \rightarrow \infty$ .*

**Example 3.1.** In (22), we gave the refactored  $\sigma$ -value operator  $S_\sigma$  associated with the expected value function factorization of a finite state Markov decision process. Recall that we set  $\mathcal{V}$  and  $\mathbb{V}$  to be the set of functions from  $\mathbf{X}$  to  $\mathbb{R}$ . So Assumption 3.1 holds based on the analysis of Section 2.2. It is straightforward to check that this choice of  $S_\sigma$  satisfies (29) when  $\alpha$  is set to the discount factor  $\beta$  and  $\kappa \equiv 1$ . Since the operators  $W_0$  and  $W_1$  associated with this plan factorization are both monotone (see (12)), the conclusions of Theorems 3.3, 3.5 and 3.6 are all valid.

**Example 3.2.** The refactored  $\sigma$ -value operator  $S_\sigma$  corresponding to the plan factorization used in the optimal stopping problem from Section 2.4.4 was given in (23). Since both  $\mathcal{V}$  and  $\mathbb{V}$  are the set of bounded continuous functions on  $\mathbf{X}$ , Assumption 3.1 has been verified in Section 2.4.4. Moreover, it is straightforward to show that  $S_\sigma$  satisfies (29) when  $\alpha := \beta$  and  $\kappa \equiv 1$ . Hence, for this plan factorization of the optimal stopping problem, the conclusions of Theorems 3.3–3.6 hold.

#### 4. APPLICATIONS

In this section we connect the theory presented above with applications. In each case, we use transformations of the Bellman equation that reduce dimensionality and enhance computational efficiency. Throughout, we use  $\mathbb{E}_{y|x}$  to denote the expectation of  $y$  conditional on  $x$ . Code that replicates all of our numerical results can be found at [https://github.com/jstac/dp\\_deconstructed\\_code](https://github.com/jstac/dp_deconstructed_code).

**4.1. Consumer Bankruptcy.** Livshits et al. (2007) analyze consumer bankruptcy rules by building a dynamic model of household choice with earnings and expense uncertainty. In the model (slightly streamlined), a household's time  $t$  income is  $z_t\eta_t$ , where  $z_t$  and  $\eta_t$  are the persistent and transitory components of productivity. Households face an expense shock  $\kappa_t \geq 0$ . While  $\{z_t\}$  is Markov,  $\{\eta_t\}$  and  $\{\kappa_t\}$  are IID. All are mutually independent. The Bellman equation takes the form

$$v(i, d, z, \eta, \kappa) = \max_{c, d', i'} [u(c) + \beta \mathbb{E}_{z', \eta', \kappa' | z} v(i', d', z', \eta', \kappa')] , \quad (32)$$

where  $d$  is debt,  $c$  is current consumption,  $\beta$  is a discount factor, and  $u$  is one-period utility. Primes denote next period values. The state  $i$  indicates repayment status and lies in  $\{R, B, E\}$ , where  $R$ ,  $B$  and  $E$  correspond to repayment, bankruptcy and default on current expenses. With  $q$  and  $\bar{r}$  indicating the interest rate for a household in states  $R$  and  $E$  respectively, and  $\gamma$  parameterizing compulsory repayments out of current income, the constraints can be expressed as follows: All variables are nonnegative and, in addition,

- C1. if  $i = R$ , then  $c + d + \kappa = z\eta + q(d', z)d'$  and  $i' \in \{R, B\}$ .
- C2. If  $i = B$ , then  $c = (1 - \gamma)z\eta$ ,  $d' = 0$  and  $i' \in \{R, E\}$ .
- C3. If  $i = E$ , then  $c = (1 - \gamma)z\eta$ ,  $d' = (\kappa - \gamma z\eta)(1 + \bar{r})$  and  $i' \in \{R, B\}$ .

A full interpretation of these constraints and the decision problem of the household can be found on p. 407 of Livshits et al. (2007). As in that paper, to implement the model computationally, we assume that all states are discrete. In particular,  $\eta$ ,  $\kappa$  and  $z$  take values in finite sets, and the choice of debt level  $d$  is restricted to a finite grid with maximum  $\bar{d}$ . The model then becomes a special case of the finite state Markov decision process described in Example 2.1. The one-period reward is  $u(c)$  and the feasible correspondence  $\Gamma$  is defined by C1–C3 above. The state is  $x = (i, d, z, \eta, \kappa)$  and the state space  $\mathbf{X}$  is the Cartesian product of  $\{R, B, E\}$ , for  $i$ , and the grids for the remaining variables. The action is  $a = (i', d', c)$  and the sets of candidate value functions are  $\mathcal{V} = \mathbb{V} = \mathbb{R}^{\mathbf{X}}$ .

A fast algorithm for solving this model is important because it allows for effective exploration of the parameter space during the estimation step (by solving the model many times at different parameterizations). We can greatly accelerate value function iteration for this model by implementing a suitable plan factorization and then iterating with the refactored Bellman operator  $S$  instead of the original Bellman operator  $T$ . To see why this can deliver acceleration, recall from Figure 1 and the definition of  $S$  in (9) that one iteration of  $S$  is a cycle starting from  $\mathbb{G} = W_0\mathbb{V}$ , while  $T$  is a cycle starting from  $\mathbb{V}$ . If functions in  $\mathbb{G}$  are simpler than functions in  $\mathbb{V}$  then the former will, in general, be faster than the latter. Similar comments apply when we consider applying  $S_\sigma$  instead of  $T_\sigma$ , as occurs during policy function iteration.

For example, suppose that we adopt the plan factorization associated with the expected value function transformation from Section 2.4.3. In view of the discussion in Example 3.1, we know that refactored value function iteration and refactored policy function iteration are convergent under this plan factorization. Moreover, since,  $\mathbb{G} = W_0\mathbb{V}$ , a typical element of  $\mathbb{G}$  is a function  $g = W_0 v$ , which, in view of the definition of the expected value function in (12), has the form  $g(x, a) = \sum_{x' \in \mathcal{X}} v(x') p(x, a, x')$ . In the setting of (32), this becomes  $g(i', d', z) = \mathbb{E}_{z', \eta', \kappa' | z} v(i', d', z', \eta', \kappa')$ . Significantly, while the standard value function has five arguments, this refactored value function has only three.

To test the expected speed gains, we perform two groups of numerical experiments. In each case we compare traditional value function iteration (VFI, iterating with  $T$ ) with refactored value function iteration (RVFI, iterating with  $S$ ). Regarding the primitives, we set  $\gamma = 0.355$ ,  $\bar{r} = 0.2$  and  $u(c) = c^{1-\sigma}/(1-\sigma)$  with  $\sigma = 2.0$ . We assume that  $\{\eta_t\} \stackrel{\text{iid}}{\sim} N(0, 0.043)$  and  $\{z_t\}$  follows  $\log z_{t+1} = \rho \log z_t + \varepsilon_{t+1}$  with  $\{\varepsilon_t\} \stackrel{\text{iid}}{\sim} N(0, \delta^2)$ . Both are discretized via the method of Tauchen (1986). We discretize  $\{\kappa_t\}$  and  $\{d_t\}$  in equidistant steps, where  $\{\kappa_t\}$  is uniformly distributed on  $[0, 2]$  and the grid points for  $d$  lie in  $[0, 10]$ . In Group-1, we set  $\rho = 0.99$  and  $\delta^2 = 0.007$  and compare the time taken of VFI and RVFI for different levels of grid points and  $\beta$  values. In Group-2, we set  $\beta = 0.98$  and continue the analysis by comparing RVFI and VFI across different grid sizes and  $(\rho, \delta)$  values. The parameter values are broadly in line with Livshits et al. (2007). Each scenario, we terminate iteration when distance between successive iterates falls below  $10^{-4}$  under the supremum norm. (The code uses Python with Numba on a workstation with a 2.9 GHz Intel Core i7 and 32GB RAM.)

The results are shown in Tables 1–2. In particular, *ratio* therein represents the ratio of time taken by VFI to that by RVFI. In both groups of experiments, reduced dimensionality significantly enhances computational efficiency. The speed improvement of RVFI over VFI reaches two orders of magnitude under a relatively sparse grid and becomes

TABLE 1. Time taken in seconds: Group-1 experiments

grid size for $(d, z, \eta, \kappa)$	method	$\beta = 0.94$	$\beta = 0.95$	$\beta = 0.96$	$\beta = 0.97$	$\beta = 0.98$
(10, 10, 10, 10)	VFI	20.75	24.17	30.13	40.07	59.77
	RVFI	0.88	0.89	1.07	1.23	1.49
	Ratio	<b>23.58</b>	<b>27.16</b>	<b>28.16</b>	<b>32.58</b>	<b>40.11</b>
(12, 12, 12, 12)	VFI	92.78	110.06	138.15	184.59	277.13
	RVFI	1.46	1.56	1.81	2.23	3.00
	Ratio	<b>63.55</b>	<b>70.55</b>	<b>76.33</b>	<b>82.78</b>	<b>92.38</b>
(14, 14, 14, 14)	VFI	321.53	387.00	484.94	648.58	978.66
	RVFI	2.55	2.91	3.49	4.65	6.40
	Ratio	<b>126.09</b>	<b>132.99</b>	<b>138.95</b>	<b>139.48</b>	<b>152.92</b>
(16, 16, 16, 16)	VFI	1445.53	1738.22	2175.40	2904.84	4381.13
	RVFI	4.92	5.75	7.12	9.45	13.65
	Ratio	<b>293.81</b>	<b>302.30</b>	<b>305.53</b>	<b>307.39</b>	<b>320.96</b>
(18, 18, 18, 18)	VFI	2412.84	2889.44	3614.84	4865.62	7266.21
	RVFI	10.94	12.92	16.14	21.99	33.06
	Ratio	<b>220.55</b>	<b>223.64</b>	<b>223.97</b>	<b>221.27</b>	<b>219.79</b>
(20, 20, 20, 20)	VFI	5591.37	6737.69	8420.48	11355.90	17020.04
	RVFI	19.78	23.80	31.16	41.00	62.09
	Ratio	<b>282.68</b>	<b>283.10</b>	<b>270.23</b>	<b>276.97</b>	<b>274.12</b>

larger as the number of grid points increase—which is precisely when computational efficiency is required. For example, when  $(\rho, \delta) = (0.995, 0.2)$  and there are 16 grid points for each state variable, RVFI is 525 times faster than VFI.

**4.2. Optimal Stopping.** Recall the plan factorization used in the optimal stopping problem from Section 2.4.4, corresponding to (14). We showed in Example 3.2 that the conclusions of Theorems 3.3–3.6 hold, so we can freely use either the traditional Bellman operator or the refactored version to compute the optimal policy. The best choice depends on relative numerical efficiency of the two operators.

One setting where the refactored Bellman operator is always more numerically efficient is when the state can be decomposed as

$$x_t = (y_t, z_t) \in Y \times Z \subset \mathbb{R}^\ell \times \mathbb{R}^k, \quad (33)$$

where  $(y_{t+1}, z_{t+1})$  and  $y_t$  are independent given  $z_t$ .

TABLE 2. Time taken in seconds: Group-2 experiments

grid size for $(d, z, \eta, \kappa)$	method	$\rho = 0.96$	$\rho = 0.97$	$\rho = 0.98$	$\rho = 0.995$	
(10, 10, 10, 10)	VFI	60.83	59.66	57.69	66.70	$\delta^2 = 0.01$
	RVFI	1.49	1.43	1.43	1.49	
	Ratio	<b>40.83</b>	<b>41.72</b>	<b>40.34</b>	<b>44.77</b>	
	VFI	68.43	62.49	59.73	75.71	$\delta^2 = 0.04$
	RVFI	1.65	1.48	1.44	1.39	
	Ratio	<b>41.47</b>	<b>42.22</b>	<b>41.48</b>	<b>54.47</b>	
(12, 12, 12, 12)	VFI	255.11	274.48	268.12	327.55	$\delta^2 = 0.01$
	RVFI	2.98	2.89	2.96	3.36	
	Ratio	<b>85.61</b>	<b>94.98</b>	<b>90.58</b>	<b>97.49</b>	
	VFI	282.59	272.91	266.27	293.99	$\delta^2 = 0.04$
	RVFI	3.46	2.96	2.92	3.05	
	Ratio	<b>81.67</b>	<b>92.20</b>	<b>91.19</b>	<b>96.39</b>	
(14, 14, 14, 14)	VFI	881.84	993.94	934.95	1012.90	$\delta^2 = 0.01$
	RVFI	6.25	6.53	6.18	6.66	
	Ratio	<b>141.09</b>	<b>152.21</b>	<b>151.29</b>	<b>152.09</b>	
	VFI	1045.88	892.71	928.10	984.25	$\delta^2 = 0.04$
	RVFI	7.47	6.24	6.18	6.64	
	Ratio	<b>140.01</b>	<b>143.06</b>	<b>150.18</b>	<b>148.23</b>	
(16, 16, 16, 16)	VFI	2739.10	2722.26	3113.68	7827.82	$\delta^2 = 0.01$
	RVFI	12.81	12.75	14.50	15.03	
	Ratio	<b>213.83</b>	<b>213.51</b>	<b>214.74</b>	<b>520.81</b>	
	VFI	2843.06	2683.64	2984.70	7190.32	$\delta^2 = 0.04$
	RVFI	13.19	12.60	13.63	13.69	
	Ratio	<b>215.55</b>	<b>212.99</b>	<b>218.98</b>	<b>525.22</b>	

The refactored Bellman operator from (14) then reduces to

$$(Sg)(z) = \int \max \{r(y', z'), c(y', z') + \beta g(z')\} P(z, d(y', z')),$$

while the standard Bellman operator is

$$(Tv)(y, z) = \max \left\{ r(y, z), c(y, z) + \beta \int v(y', z') P(z, d(y', z')) \right\}.$$

The key difference is that the functions  $v \in \mathbb{V}$ , on which  $T$  acts, depend on both  $z$  and  $y$ , whereas the functions  $g \in \mathbb{G}$  on which  $S$  acts depend only on  $z$ . In other words, the refactored Bellman operator acts on functions defined over a lower dimension space than the original state space  $\mathbf{X} = \mathbf{Y} \times \mathbf{Z}$ . Hence, when the conditional independence

assumption in (33) is valid, the curse of dimensionality is mitigated by working with  $S$  rather than  $T$ . Below we quantify the difference, after giving some examples of settings where (33) holds.

**Example 4.1.** Consider the problem of pricing a real or financial option (see, e.g., Dixit and Pindyck (1994) or Kellogg (2014)). Let  $p_t$  be the current price of the asset and let  $\{s_t\}$  be a Markov process affecting asset price via  $p_t = f(s_t, \varepsilon_t)$ , where  $\{\varepsilon_t\}$  is an IID innovation process independent of  $\{s_t\}$ . Each period, the agent decides whether to exercise the option now (i.e., purchase the asset at the strike price) or wait and reconsider next period. Note that  $(p_{t+1}, s_{t+1})$  and  $p_t$  are independent given  $s_t$ . Hence, this problem fits into the framework of (33) if we take  $y_t = p_t$  and  $z_t = s_t$ .

**Example 4.2.** Consider a model of job search (see, e.g., McCall (1970)), in which a worker receives current wage offer  $w_t$  and chooses to either accept and work permanently at that wage, or reject the offer, receive unemployment compensation  $\eta_t$  and reconsider next period. The worker aims to find an optimal decision rule that yields maximum lifetime rewards. A typical formulation for the wage and compensation processes is that both are functions of an exogenous Markov process  $\{s_t\}$ , with  $w_t = f(s_t, \varepsilon_t)$  and  $\eta_t = \ell(s_t, \xi_t)$  where  $f$  and  $\ell$  are nonnegative continuous functions and  $\{\varepsilon_t\}$  and  $\{\xi_t\}$  are independent, IID innovation processes (see, e.g., Low et al. (2010), Bagger et al. (2014)). The state space for the job search problem is typically set to  $(w_t, \eta_t, s_t)$ , where  $s_t$  is included because it can be used to forecast future draws of  $w_t$  and  $\eta_t$ . The conditional independence assumption in (33) holds if we take  $y_t = (w_t, \eta_t)$  and  $z_t = s_t$ .

Let us now compare the time complexity of VFI and RVFI, based on iterating  $T$  and  $S$  respectively.

4.2.1. *Finite State Case.* Let  $\mathbf{Y} = \times_{i=1}^{\ell} \mathbf{Y}^i$  and  $\mathbf{Z} = \times_{j=1}^k \mathbf{Z}^j$ , where  $\mathbf{Y}^i$  and  $\mathbf{Z}^j$  are subsets of  $\mathbb{R}$ . Each  $\mathbf{Y}^i$  (resp.,  $\mathbf{Z}^j$ ) is represented by a grid of  $L_i$  (resp.,  $K_j$ ) points. Integration operations in both VFI and RVFI are replaced by summations. We use  $\hat{P}$  and  $\hat{F}$  to denote the probability transition matrices for VFI and RVFI respectively.

Let  $L := \prod_{i=1}^{\ell} L_i$  and  $K := \prod_{j=1}^k K_j$  with  $L = 1$  for  $\ell = 0$ . Let  $k > 0$ . There are  $LK$  grid points on  $\mathbf{X} = \mathbf{Y} \times \mathbf{Z}$  and  $K$  grid points on  $\mathbf{Z}$ . The matrix  $\hat{P}$  is  $(LK) \times (LK)$  and  $\hat{F}$  is  $K \times (LK)$ . VFI and RVFI are implemented by the operators  $\hat{T}$  and  $\hat{S}$  defined respectively by

$$\hat{T}\vec{v} := \vec{r} \vee (\vec{c} + \beta \hat{P}\vec{v}) \quad \text{and} \quad \hat{S}\vec{g}_z := \hat{F}[\vec{r} \vee (\vec{c} + \beta \vec{g})].$$

TABLE 3. Time complexity: VFI v.s RVFI

Complexity	VFI: 1-loop	RVFI: 1-loop	VFI: $n$ -loop	RVFI: $n$ -loop
FS	$\mathcal{O}(L^2 K^2)$	$\mathcal{O}(LK^2)$	$\mathcal{O}(nL^2 K^2)$	$\mathcal{O}(nLK^2)$
CS	$\mathcal{O}(NLK \log(LK))$	$\mathcal{O}(NK \log(K))$	$\mathcal{O}(nNLK \log(LK))$	$\mathcal{O}(nNK \log(K))$

Here  $\vec{q}$  represents a column vector with  $i$ -th element equal to  $q(y_i, z_i)$ , where  $(y_i, z_i)$  is the  $i$ -th element of the list of grid points on  $\mathbf{Y} \times \mathbf{Z}$ .  $\vec{q}_z$  denotes the column vector with the  $j$ -th element equal to  $q(z_j)$ , where  $z_j$  is the  $j$ -th element of the list of grid points on  $\mathbf{Z}$ . The vectors  $\vec{v}$ ,  $\vec{r}$ ,  $\vec{c}$  and  $\vec{g}$  are  $(LK) \times 1$ , while  $\vec{g}_z$  is  $K \times 1$ . Moreover,  $\vee$  denotes taking maximum.

4.2.2. *Continuous State Case.* We use the same number of grid points as before, but now for continuous state function approximation rather than discretization. In particular, we replace the discrete state summation with Monte Carlo integration. Assume that the transition law of the state process follows

$$y_{t+1} = f_1(z_t, \varepsilon_{t+1}), \quad z_{t+1} = f_2(z_t, \varepsilon_{t+1}), \quad \{\varepsilon_t\} \stackrel{\text{iid}}{\sim} \Phi.$$

After drawing  $N$  Monte Carlo samples  $U_1, \dots, U_N \stackrel{\text{iid}}{\sim} \Phi$ , RVFI and VFI are implemented by

$$\hat{S}g(z) := \frac{1}{N} \sum_{i=1}^N \max \{ r(f_1(z, U_i), f_2(z, U_i)), c(f_1(z, U_i), f_2(z, U_i)) + \beta \phi\langle g \rangle(f_2(z, U_i)) \}$$

$$\text{and } \hat{T}v(y, z) := \max \left\{ r(y, z), c(y, z) + \beta \frac{1}{N} \sum_{i=1}^N \psi\langle v \rangle(f_1(z, U_i), f_2(z, U_i)) \right\}.$$

Here  $g = \{g(z)\}$  with  $z$  in the set of grid points on  $\mathbf{Z}$ , and  $v = \{v(y, z)\}$  with  $(y, z)$  in the set of grid points on  $\mathbf{Y} \times \mathbf{Z}$ . Moreover,  $\phi\langle \cdot \rangle$  and  $\psi\langle \cdot \rangle$  are interpolating functions for RVFI and VFI respectively. For example,  $\phi\langle g \rangle(z)$  interpolates the vector  $g$  to obtain a function  $\phi\langle g \rangle$  and then evaluates it at  $z$ .

4.2.3. *Time Complexity.* Table 3 provides the time complexity of RVFI and VFI, estimated by counting the number of floating point operations. Each such operation is assumed to have complexity  $\mathcal{O}(1)$ , so is function evaluation associated with the model primitives. Moreover, for continuous state case, binary search is used when we evaluate the interpolating function at a given point, and our results hold for linear, quadratic, cubic, and  $k$ -nearest neighbors interpolations.

For both finite state (FS) and continuous state (CS) cases, RVFI provides better performance than VFI. For FS, RVFI is more efficient than VFI by order  $\mathcal{O}(L)$ , while for CS,

RVFI is so by order  $\mathcal{O}(L \log(LK)/\log(K))$ . For example, if we have 100 grid points in each dimension, in the FS case, evaluating a given number of loops will take around  $100^\ell$  times longer via VFI than via RVFI, after adjusting for order approximations. See the Appendix (Section 6.3) for proof of Table 3 results.

**4.3. Robust Control.** In the robust control problem studied by Bidder and Smith (2012) (see Hansen and Sargent (2008) for systematic discussion), an agent is endowed with a benchmark model (of the system transition probabilities) but fears that it is misspecified. While making decisions, she considers alternative distributions that are distorted versions of the distribution implied by the benchmark, and balances the cost that an implicit misspecification would cause against the plausibility of misspecification. The state-action aggregator is given by

$$H(s, \varepsilon, u, v) = r(s, \varepsilon, u) - \beta\theta \log \left[ \mathbb{E}_{s', \varepsilon' | s} \exp \left( -\frac{v(s', \varepsilon')}{\theta} \right) \right] \quad \text{with} \quad s' = f(s, u, \varepsilon'),$$

where  $s_t$  is an endogenous state component and  $u_t$  is a vector of controls. The full state is  $x_t = (s_t, \varepsilon_t)$ , where  $\{\varepsilon_t\}$  is an IID innovation process. Rewards at time  $t$  are  $r(s_t, \varepsilon_t, u_t)$ . Let  $\theta > 0$  be the penalty parameter that controls the degree of robustness. The Bellman operator is

$$Tv(s, \varepsilon) = \max_u \left\{ r(s, \varepsilon, u) - \beta\theta \log \left[ \mathbb{E}_{s', \varepsilon' | s} \exp \left( -\frac{v(s', \varepsilon')}{\theta} \right) \right] \right\}.$$

Consider the plan factorization  $(W_0, W_1)$  defined by

$$W_0 v(s) := -\theta \log \left\{ \mathbb{E}_{s', \varepsilon' | s} \exp \left[ -\frac{v(s', \varepsilon')}{\theta} \right] \right\} \quad \text{and} \quad W_1 g(s, \varepsilon, u) := r(s, \varepsilon, u) + \beta g(s).$$

The refactored Bellman operator  $S = W_0 M W_1$  is then

$$Sg(s) = -\theta \log \left\{ \mathbb{E}_{s', \varepsilon' | s} \exp \left[ -\frac{\max_{u'} \{r(s', \varepsilon', u') + \beta g(s')\}}{\theta} \right] \right\}.$$

The maps  $W_0$  and  $W_1$  are monotone in the standard pointwise ordering for real valued functions, and the other assumptions of Sections 2–3 hold in the setting of Bidder and Smith (2012), so the conclusions of Theorems 3.3–3.6 hold. The details are omitted. While the value function acts on both  $s$  and  $\varepsilon$ , the refactored value function acts only on  $s$ . Hence, the refactored Bellman operator mitigates the curse of dimensionality, similar to the outcomes in Sections 4.1 and 4.2.



## 5. CONCLUSION

This paper presents a systematic treatment of the technique of transforming Bellman equations associated with discrete time dynamic programs in order to convert them into more advantageous forms. We formalized these transformations in terms of what we called plan factorizations and showed how this concept encompasses and extends existing transformations from the literature. We provided conditions related to monotonicity under which the transformations preserve the most valuable property of Bellman equations: their link to optimal value functions and optimal decisions.

The applications presented here focused on how transforming the Bellman equation can mitigate the curse of dimensionality and thereby boost computational efficiency. There are, however, other reasons to consider transformations of the Bellman equation. For example, it might be that a given transformation leads to smoother value functions. Smoother functions are easier to approximate in high dimensional spaces. It also appears that some transformations of a given Bellman equation can shift the problem from a setting of unbounded value functions to bounded ones, where characterization of optimality becomes easier. While these ideas go beyond the scope of the current study, the theory of plan factorizations presented here should serve as a solid foundation for new work along these lines.

## 6. APPENDIX

The appendix provides all remaining proofs.

**6.1. Preliminaries.** Let  $E_i$  be a nonempty set and let  $\tau_i$  be a mapping from  $E_i$  to  $E_{i+1}$  for  $i = 0, 1, 2$  with addition modulo 3 (a convention we adopt throughout this section). Consider the self-mappings

$$F_0 := \tau_2 \tau_1 \tau_0, \quad F_1 := \tau_0 \tau_2 \tau_1 \quad \text{and} \quad F_2 := \tau_1 \tau_0 \tau_2$$

on  $E_0$ ,  $E_1$  and  $E_2$  respectively. We then have

$$F_{i+1} \tau_i = \tau_i F_i \quad \text{on } E_i \text{ for } i = 0, 1, 2. \tag{34}$$

**Lemma 6.1.** *For each  $i = 0, 1, 2$ ,*

- (1) *if  $\varrho$  is a fixed point of  $F_i$  in  $E_i$ , then  $\tau_i \varrho$  is a fixed point of  $F_{i+1}$  in  $E_{i+1}$ .*
- (2)  *$F_{i+1}^n \tau_i = \tau_i F_i^n$  on  $E_i$  for all  $n \in \mathbb{N}$ .*

*Proof.* Regarding part (a), if  $\varrho$  is a fixed point of  $F_i$  in  $E_i$ , then (34) yields  $F_{i+1}\tau_i\varrho = \tau_i F_i\varrho = \tau_i\varrho$ , so  $\tau_i\varrho$  is a fixed point of  $F_{i+1}$ . Regarding part (b), fix  $i$  in  $\{0, 1, 2\}$ . By (34), the statement in (b) is true at  $n = 1$ . Let it also be true at  $n - 1$ . Then, using (34) again,  $F_{i+1}^n\tau_i = F_{i+1}^{n-1}F_{i+1}\tau_i = F_{i+1}^{n-1}\tau_i F_i = \tau_i F_i^{n-1}F_i = \tau_i F_i^n$ . We conclude that (b) holds at every  $n \in \mathbb{N}$ .  $\square$

**Lemma 6.2.** *If  $F_i$  has a unique fixed point  $\varrho_i$  in  $E_i$  for some  $i$  in  $\{0, 1, 2\}$ , then  $\tau_i\varrho_i$  is the unique fixed point of  $F_{i+1}$  in  $E_{i+1}$ .*

*Proof.* We have already proved all but uniqueness. To see that uniqueness holds, fix  $i \in \{0, 1, 2\}$  and suppose that  $F_i$  has only one fixed point in  $E_i$ , whereas  $F_{i+1}$  has two fixed points in  $E_{i+1}$ . Denote the fixed points of  $F_{i+1}$  by  $\varrho$  and  $f$ . Applying part (a) of Lemma 6.1 twice, we see that  $\tau_{i+2}\tau_{i+1}\varrho$  and  $\tau_{i+2}\tau_{i+1}f$  are both fixed points of  $F_{i+3} = F_i$ . Since  $F_i$  has only one fixed point,  $\tau_{i+2}\tau_{i+1}\varrho = \tau_{i+2}\tau_{i+1}f$ . Applying  $\tau_i$  to both sides of the last equality gives  $F_{i+1}\varrho = F_{i+1}f$ . Since  $\varrho$  and  $f$  are both fixed points of  $F_{i+1}$ , we conclude that  $\varrho = f$  and the fixed point is unique.  $\square$

**6.2. Proof of Sections 2–3 Results.** When connecting to the results in Section 6.1, we always take  $\tau_0 = W_0$  and  $\tau_1 = W_1$ . The map  $\tau_2$  and the candidate spaces  $E_0$ ,  $E_1$  and  $E_2$  will vary depending on the context.

*Proof of Proposition 2.1.* The claims are immediate from Lemma 6.1. Regarding iterations on  $S$  and  $T$ , we set  $E_0 = \mathbb{V}$ ,  $E_1 = \mathbb{G}$ ,  $E_2 = \mathbb{H}$  and  $\tau_2 = M$ . Regarding iterations on  $S_\sigma$  and  $T_\sigma$ , we fix  $\sigma \in \Sigma$  and set  $E_0 = \mathcal{V}$ ,  $E_1 = \mathcal{G}$ ,  $E_2 = \mathcal{H}$  and  $\tau_2 = M_\sigma$ .  $\square$

*Proof of Proposition 2.2.* The claims are immediate by Lemmas 6.1–6.2 (let  $E_0 = \mathbb{V}$ ,  $E_1 = \mathbb{G}$ ,  $E_2 = \mathbb{H}$  and  $\tau_2 = M$ ).  $\square$

*Proof of Proposition 2.3.* Similar to the proof of Proposition 2.2, this result is immediate from Lemmas 6.1–6.2 once we set  $E_0 = \mathcal{V}$ ,  $E_1 = \mathcal{G}$ ,  $E_2 = \mathcal{H}$  and  $\tau_2 = M_\sigma$ .  $\square$

*Proof of Theorem 3.2.* Now suppose that (a) of Theorem 3.2 holds, so that  $g^* \in \mathbb{G}$  and  $Sg^* = g^*$ . We claim that (b) holds, i.e.,  $v^* \in \mathbb{V}$  and  $Tv^* = v^*$ . Based on Lemma 6.1 (in particular, let  $E_0 = \mathbb{V}$ ,  $E_1 = \mathbb{G}$ ,  $E_2 = \mathbb{H}$  and  $\tau_2 = M$ ), we only need to show that  $v^* = MW_1g^*$ .

Pick any  $\sigma \in \Sigma$ . We have

$$v_\sigma = M_\sigma W_1 g_\sigma \leq MW_1 g_\sigma \leq MW_1 g^*,$$

where the equality is due to Assumption 2.2 and Proposition 2.3, the first inequality is due to the fact that  $M$  pointwise dominates  $M_\sigma$ , and the second follows from the definition of  $g^*$  and the monotonicity of  $M$  and  $W_1$ . As  $\sigma$  was arbitrary, this proves that  $v^* \leq MW_1g^*$ .

Regarding the reverse inequality, from  $g^* = Sg^* \in \mathbb{G}$ , there exists a  $v \in \mathbb{V}$  such that  $g^* = W_0v$ . Assumption 2.1 then implies the existence of a  $\sigma$  such that  $MW_1g^* = MW_1W_0v = M_\sigma W_1W_0v = M_\sigma W_1g^*$ . Hence,  $Sg^* = W_0MW_1g^* = W_0M_\sigma W_1g^* = S_\sigma g^*$ . We see that  $g^*$  is a fixed point of  $S_\sigma$  in  $\mathcal{G}$ . Since  $g_\sigma$  is the unique fixed point of  $S_\sigma$  in  $\mathcal{G}$  by Assumption 2.2 and Proposition 2.3, we have  $g^* = g_\sigma$ . As a result,  $MW_1g^* = M_\sigma W_1g_\sigma = v_\sigma \leq v^*$ . We have now shown that (b) holds.

The claim that (b) implies (a) is similar. In particular, based on Lemma 6.1, it suffices to show that  $g^* = W_0v^*$ . Pick any  $\sigma \in \Sigma$ . Since  $g_\sigma = W_0v_\sigma \leq W_0v^*$  by Assumption 2.2, Proposition 2.3 and the monotonicity of  $W_0$ , it follows that  $g^* \leq W_0v^*$ . Moreover, since  $v^* \in \mathbb{V}$ , Assumption 2.1 implies the existence of a  $\sigma$  such that  $Tv^* = T_\sigma v^*$ . Since in addition  $v^* = Tv^*$ , it follows that  $v^*$  is a fixed point of  $T_\sigma$  and thus  $v^* = v_\sigma$  by Assumption 2.2. We then have  $W_0v^* = W_0v_\sigma = g_\sigma \leq g^*$ . In summary, we have  $g^* = W_0v^*$  and claim (a) is now verified.

The claim that  $g^* = \hat{g}$  under (a)–(b) translates to the claim that  $g^* = W_0v^*$ , which we have already shown in the previous step.

The claim in Theorem 3.2 that at least one optimal policy exists under either (a) or (b) follows from the equivalence of (a) and (b) just established combined with Theorem 3.1.

Finally, suppose that (a) holds and consider the last claim in Theorem 3.2, which can be expressed succinctly as

$$v_\sigma = v^* \iff T_\sigma v^* = Tv^* \iff M_\sigma W_1g^* = MW_1g^*.$$

Regarding the first equivalence, if  $v_\sigma = v^*$ , then  $T_\sigma v^* = T_\sigma v_\sigma = v_\sigma = v^* = Tv^*$ . Conversely, if  $T_\sigma v^* = Tv^*$ , then, since  $v^* = Tv^*$  and, by Assumption 2.2,  $v_\sigma$  is the only fixed point of  $T_\sigma$  in  $\mathcal{V}$ , we have  $v_\sigma = v^*$ . The second equivalence follows, since, by the relations  $g_\sigma = W_0v_\sigma$  and  $g^* = W_0v^*$ , we have:  $T_\sigma v^* = M_\sigma W_1W_0v^* = M_\sigma W_1g^*$  and  $Tv^* = MW_1W_0v^* = MW_1g^*$ . The last claim in Theorem 3.2 is now established. This concludes the proof.  $\square$

*Proof of Theorem 3.3.* We need only show that (a) holds, since the remaining claims follow directly from the hypotheses of the theorem, claim (a), the definition of asymptotic stability (for part (b)) and Theorem 3.2 (for parts (c) and (d)).

To see that (a) holds, note that, by the stated hypotheses,  $S$  has a unique fixed point in  $\mathbb{G}$ , which we denote below by  $\bar{g}$ . Our aim is to show that  $\bar{g} = g^*$ .

First observe that, by Assumption 2.1, there is a  $\sigma \in \Sigma$  such that  $S_\sigma \bar{g} = S\bar{g} = \bar{g}$ . But, by Assumption 2.2 and Proposition 2.3,  $S_\sigma$  has exactly one fixed point, which is the refactored  $\sigma$ -value function  $g_\sigma$ . Hence  $\bar{g} = g_\sigma$ . In particular,  $\bar{g} \leq g^*$ .

To see that the reverse inequality holds, pick any  $\sigma \in \Sigma$  and note that, by the definition of  $S$ , we have  $\bar{g} = S\bar{g} \geq S_\sigma \bar{g}$ . We know that  $S_\sigma$  is monotone on  $\mathbb{G}$ , since this operator is the composition of three monotone operators ( $W_0$  and  $W_1$  by assumption and  $M_\sigma$  automatically). Hence we can iterate on the last inequality to establish  $\bar{g} \geq S_\sigma^k \bar{g}$  for all  $k \in \mathbb{N}$ . Taking limits and using the fact that the partial order is closed,  $S_\sigma$  is asymptotically stable and, as shown above,  $g_\sigma$  is the unique fixed point, we have  $\bar{g} \geq g_\sigma$ . Since  $\sigma$  was chosen arbitrarily, this yields  $\bar{g} \geq g^*$ . Part (a) of Theorem 3.3 is now verified.  $\square$

*Proof of Proposition 3.4.* As a closed subset of  $b_\kappa \mathbf{F}$  under the  $\|\cdot\|_\kappa$ -norm metric, the set  $\mathcal{G}$  is complete under the same metric and, by (29) and the Banach contraction mapping theorem, each  $S_\sigma$  is asymptotically stable on  $\mathcal{G}$ . Moreover, the pointwise partial order  $\leq$  is closed under this metric. Thus, to verify the conditions of Theorem 3.3, we need only show that  $S$  is asymptotically stable on  $\mathbb{G}$  under the same metric. To this end, we first claim that, under the stated assumptions,

$$Sg(x, a) = \sup_{\sigma \in \Sigma} S_\sigma g(x, a) \quad \text{for all } (x, a) \in \mathbf{F} \text{ and } g \in \mathbb{G}. \quad (35)$$

To see that (35) holds, fix  $g \in \mathbb{G}$ . Notice that the definition of  $M$  implies that  $MW_1 g \geq M_\sigma W_1 g$  for any  $\sigma \in \Sigma$  and hence, applying  $W_0$  to both sides and using monotonicity yield  $Sg \geq S_\sigma g$  for all  $\sigma \in \Sigma$ . Moreover, by the definition of  $\mathbb{G}$  and Assumption 2.1, there exists a  $g$ -greedy policy  $\sigma^* \in \Sigma$  such that  $MW_1 g = M_{\sigma^*} W_1 g$ , and applying  $W_0$  to both sides again gives  $Sg = S_{\sigma^*} g$ . Hence (35) is valid.

Now fix  $g, g' \in \mathbb{G}$  and  $(x, a) \in \mathbf{F}$ . By (35) and the contraction condition on  $S_\sigma$  in (29), we have

$$\begin{aligned} |Sg(x, a) - Sg'(x, a)| &= \left| \sup_{\sigma \in \Sigma} S_\sigma g(x, a) - \sup_{\sigma \in \Sigma} S_\sigma g'(x, a) \right| \leq \sup_{\sigma \in \Sigma} |S_\sigma g(x, a) - S_\sigma g'(x, a)|, \\ \therefore \frac{|Sg(x, a) - Sg'(x, a)|}{\kappa(x, a)} &\leq \sup_{\sigma \in \Sigma} \|S_\sigma g - S_\sigma g'\|_\kappa \leq \alpha \|g - g'\|_\kappa. \end{aligned}$$

Taking supremum gives  $\|Sg - Sg'\|_\kappa \leq \alpha \|g - g'\|_\kappa$ , so  $S$  is a contraction on  $\mathbb{G}$  of modulus  $\alpha$ . Since in addition  $\mathbb{G}$  is a closed subset of  $b_\kappa \mathbf{F}$  under the  $\|\cdot\|_\kappa$ -norm metric, it is a Banach space under the same norm. The Banach contraction mapping theorem then

implies that  $S$  is asymptotically stable on  $\mathbb{G}$ . Since all the conditions of Theorem 3.3 are verified, its conclusions follow.  $\square$

*Proof of Theorem 3.5.* Without loss of generality, we assume that  $m_k \equiv m$  for all  $k \in \mathbb{N}_0$ . When  $k = 0$ , by definition,

$$\Sigma_0^v = \{\sigma \in \Sigma: M_\sigma W_1 W_0 v_0 = M W_1 W_0 v_0\} = \{\sigma \in \Sigma: M_\sigma W_1 g_0 = M W_1 g_0\} = \Sigma_0^g,$$

i.e., the sets of  $v_0$ -greedy and  $g_0$ -greedy policies are identical. This in turn implies that a policy  $\sigma_0 \in \Sigma$  is  $v_0$ -greedy if and only if it is  $g_0$ -greedy. Moreover, by Proposition 2.1,

$$\begin{aligned} g_1 &= S_{\sigma_0}^m g_0 = W_0 T_{\sigma_0}^{m-1} M_{\sigma_0} W_1 g_0 \\ &= W_0 T_{\sigma_0}^{m-1} M_{\sigma_0} W_1 W_0 v_0 = W_0 T_{\sigma_0}^{m-1} T_{\sigma_0} v_0 = W_0 T_{\sigma_0}^m v_0 = W_0 v_1. \end{aligned}$$

The related claims are verified for  $k = 0$ . Suppose these claims hold for arbitrary  $k$ . It remains to show that they hold for  $k + 1$ . By the induction hypothesis,  $\Sigma_k^v = \Sigma_k^g$ , a policy  $\sigma_k \in \Sigma$  is  $v_k$ -greedy if and only if it is  $g_k$ -greedy, and  $g_{k+1} = W_0 v_{k+1}$ . By definition of the greedy policy,

$$\begin{aligned} \Sigma_{k+1}^v &= \{\sigma \in \Sigma: M_\sigma W_1 W_0 v_{k+1} = M W_1 W_0 v_{k+1}\} \\ &= \{\sigma \in \Sigma: M_\sigma W_1 g_{k+1} = M W_1 g_{k+1}\} = \Sigma_{k+1}^g, \end{aligned}$$

i.e., a policy  $\sigma_{k+1} \in \Sigma$  is  $v_{k+1}$ -greedy if and only if it is  $g_{k+1}$ -greedy. Moreover, by Proposition 2.1,

$$\begin{aligned} g_{k+2} &= S_{\sigma_{k+1}}^m g_{k+1} = W_0 T_{\sigma_{k+1}}^{m-1} M_{\sigma_{k+1}} W_1 g_{k+1} \\ &= W_0 T_{\sigma_{k+1}}^{m-1} M_{\sigma_{k+1}} W_1 W_0 v_{k+1} = W_0 T_{\sigma_{k+1}}^m v_{k+1} = W_0 v_{k+2}. \end{aligned}$$

Hence, the related claims of the proposition hold for  $k + 1$ , completing the proof.  $\square$

*Proof of Theorem 3.6.* Without loss of generality, we assume  $m_k \equiv m$  for all  $k \in \mathbb{N}_0$ . Let  $\{\zeta_k\}$  be defined by  $\zeta_0 := g_0$  and  $\zeta_k := S \zeta_{k-1}$ . We show by induction that

$$\zeta_k \leq g_k \leq g^* \quad \text{and} \quad g_k \leq S g_k \quad \text{for all } k \in \mathbb{N}_0. \quad (36)$$

Note that  $g_0 \leq S g_0$  by assumption. Since  $S$  is monotone, this implies that  $g_0 \leq S^t g_0$  for all  $t \in \mathbb{N}_0$ . Letting  $t \rightarrow \infty$ , Theorem 3.3 implies that  $g_0 \leq g^*$ . Since in addition  $\zeta_0 = g_0$ , (36) is satisfied for  $k = 0$ . Suppose this claim holds for arbitrary  $k \in \mathbb{N}_0$ . Then, since  $S$  and  $S_{\sigma_k}$  are monotone,

$$g_{k+1} = S_{\sigma_k}^m g_k \leq S S_{\sigma_k}^{m-1} g_k \leq S S_{\sigma_k}^{m-1} S g_k = S S_{\sigma_k}^m g_k = S g_{k+1},$$

where the first inequality holds by the definition of  $S$ , the second inequality holds since  $g_k \leq S g_k$  (the induction hypothesis), and the second equality is due to  $\sigma_k \in \Sigma_k^g$ . Hence,  $g_{k+1} \leq S^t g_{k+1}$  for all  $t \in \mathbb{N}_0$ . Letting  $t \rightarrow \infty$  yields  $g_{k+1} \leq g^*$ .

Similarly, since  $g_k \leq Sg_k = S_{\sigma_k}g_k$ , we have  $g_k \leq S_{\sigma_k}^{m-1}g_k$  and thus

$$\zeta_{k+1} = S\zeta_k \leq Sg_k = S_{\sigma_k}g_k \leq S_{\sigma_k}^m g_k = g_{k+1}.$$

Hence, (36) holds by induction. The monotonicity of  $S_{\sigma_k}$  implies that  $g_k \leq Sg_k = S_{\sigma_k}g_k \leq \dots \leq S_{\sigma_k}^m g_k = g_{k+1}$  for all  $k \in \mathbb{N}_0$ . Hence,  $\{g_k\}$  is increasing in  $k$ . Moreover, since  $\zeta_k \rightarrow g^*$  by Theorem 3.3, (36) implies that  $g_k \rightarrow g^*$  as  $k \rightarrow \infty$ .  $\square$

**6.3. Proof of Section 4.2 Results.** To prove the results of Table 3, recall that floating point operations are any elementary actions (e.g.,  $+$ ,  $\times$ ,  $\max$ ,  $\min$ ) on or assignments with floating point numbers. If  $f$  and  $h$  are scalar functions on  $\mathbb{R}^n$ , we write  $f(x) = \mathcal{O}(h(x))$  whenever there exist  $C, M > 0$  such that  $\|x\| \geq M$  implies  $|f(x)| \leq C|h(x)|$ , where  $\|\cdot\|$  is the sup norm. We introduce some elementary facts on time complexity:

- (a) The multiplication of an  $m \times n$  matrix and an  $n \times p$  matrix has complexity  $\mathcal{O}(mnp)$ . (See, e.g., Section 2.5.4 of Skiena (2008).)
- (b) The binary search algorithm finds the index of an element in a given sorted array of length  $n$  in  $\mathcal{O}(\log(n))$  time. (See, e.g., Section 4.9 of Skiena (2008).)

For the finite state (FS) case, time complexity of 1-loop VFI reduces to the complexity of matrix multiplication  $\hat{P}\vec{v}$ , which is of order  $\mathcal{O}(L^2K^2)$  based on the shape of  $\hat{P}$  and  $\vec{v}$  and fact (a) above. Similarly, time complexity of 1-loop RVFI has complexity  $\mathcal{O}(LK^2)$ . Time complexity of  $n$ -loop algorithms is scaled up by  $\mathcal{O}(n)$ .

For the continuous state (CS) case, let  $\mathcal{O}(\phi)$  and  $\mathcal{O}(\psi)$  denote respectively the complexity of single point evaluation of the interpolating functions  $\phi$  and  $\psi$ . Counting the floating point operations associated with all grid points inside the inner loops shows that the one step complexities of VFI and RVFI are  $\mathcal{O}(NLK)\mathcal{O}(\psi)$  and  $\mathcal{O}(NK)\mathcal{O}(\phi)$ , respectively. Since binary search function evaluation is adopted for the indicated function interpolation mechanisms, and in particular, since evaluating  $\psi$  at a given point uses binary search  $\ell + k$  times, based on fact (b) above, we have

$$\begin{aligned} \mathcal{O}(\psi) &= \mathcal{O}\left(\sum_{i=1}^{\ell} \log(L_i) + \sum_{j=1}^k \log(K_j)\right) \\ &= \mathcal{O}\left(\log\left(\prod_{i=1}^{\ell} L_i\right) + \log\left(\prod_{j=1}^k K_j\right)\right) = \mathcal{O}(\log(LK)). \end{aligned}$$

Similarly, we can show that  $\mathcal{O}(\phi) = \mathcal{O}(\log(K))$ . Combining these results, we see that the claims of the CS case hold, concluding our proof of Table 3 results.

**6.4. Counterexamples.** In this section, we provide counterexamples showing that monotonicity of  $W_0$  and  $W_1$  cannot be dropped from Theorem 3.2. Recall that  $\mathbb{E}_{y|x}$  denotes the expectation of  $y$  conditional on  $x$ .

**6.4.1. Counterexample 1.** Here we exhibit a dynamic program and value transformation under which  $Tv^* = v^*$ ,  $Sg^* \neq g^*$  and  $g^* \neq \hat{g}$ . The example involves risk sensitive preferences.

Let  $\mathbf{X} := \{1, 2\}$ ,  $\mathbf{A} := \{0, 1\}$ ,  $\mathcal{V} = \mathbb{V} := \mathbb{R}^{\mathbf{X}}$  and  $\Sigma := \{\sigma_1, \sigma_2, \sigma_3, \sigma_4\}$ , where the policy functions are

$x$	$\sigma_1(x)$	$\sigma_2(x)$	$\sigma_3(x)$	$\sigma_4(x)$
1	0	0	1	1
2	0	1	1	0

In state  $x$ , choosing action  $a$  gives the agent an immediate reward  $x - a$ . If  $a = 0$ , then the next period state  $x' = 1$ , while if  $a = 1$ , the next period state  $x' = 2$ . Let the state-action aggregator  $H$  take the risk-sensitive form

$$H(x, a, v) = x - a - \frac{\beta}{\gamma} \log \mathbb{E}_{x'|a} \exp[-\gamma v(x')].$$

Then the Bellman operator  $T$  is

$$Tv(x) = \max_{a \in \{0,1\}} \left\{ x - a - \frac{\beta}{\gamma} \log \mathbb{E}_{x'|a} \exp[-\gamma v(x')] \right\} = \max \{x + \beta v(1), x - 1 + \beta v(2)\}.$$

Consider the plan factorization  $(W_0, W_1)$  given by

$$W_0 v(a) := \mathbb{E}_{x'|a} \exp[-\gamma v(x')] \quad \text{and} \quad W_1 g(x, a) := x - a - \frac{\beta}{\gamma} \log g(a).$$

The refactored Bellman operator is then

$$Sg(a) = \mathbb{E}_{x'|a} \exp \left( -\gamma \max_{a' \in \{0,1\}} \left\{ x' - a' - \frac{\beta}{\gamma} \log g(a') \right\} \right).$$

Note that neither  $W_1$  nor  $W_0$  is monotone. In the following, we assume that  $\beta \in (0, 1)$ .

**Lemma 6.3.** *The  $\sigma$ -value functions are as follows:*

$x$	$v_{\sigma_1}(x)$	$v_{\sigma_2}(x)$	$v_{\sigma_3}(x)$	$v_{\sigma_4}(x)$
1	$1/(1 - \beta)$	$1/(1 - \beta)$	$\beta/(1 - \beta)$	$2\beta/(1 - \beta^2)$
2	$(2 - \beta)/(1 - \beta)$	$1/(1 - \beta)$	$1/(1 - \beta)$	$2/(1 - \beta^2)$

*Proof.* Regarding  $v_{\sigma_1}$ , by definition,

$$v_{\sigma_1}(1) = 1 - \sigma_1(1) - \frac{\beta}{\gamma} \log \mathbb{E}_{x'|\sigma_1(1)} \exp[-\gamma v_{\sigma_1}(x')] = 1 + \beta v_{\sigma_1}(1)$$

$$\text{and } v_{\sigma_1}(2) = 2 - \sigma_1(2) - \frac{\beta}{\gamma} \log \mathbb{E}_{x'|\sigma_1(2)} \exp[-\gamma v_{\sigma_1}(x')] = 2 + \beta v_{\sigma_1}(1).$$

Hence,  $v_{\sigma_1}(1) = 1/(1 - \beta)$  and  $v_{\sigma_1}(2) = (2 - \beta)/(1 - \beta)$ . The remaining proofs are similar.  $\square$

Based on Lemma 6.3, the value function  $v^*$  is given by

$$v^*(1) = \max_{\sigma_i} v_{\sigma_i}(1) = v_{\sigma_1}(1) = v_{\sigma_2}(1) = 1/(1 - \beta)$$

$$\text{and } v^*(2) = \max_{\sigma_i} v_{\sigma_i}(2) = v_{\sigma_1}(2) = (2 - \beta)/(1 - \beta).$$

Moreover, we have  $v^* = Tv^*$ , since

$$Tv^*(1) = \max\{1 + \beta v^*(1), \beta v^*(2)\} = 1/(1 - \beta) = v^*(1)$$

$$\text{and } Tv^*(2) = \max\{2 + \beta v^*(1), 1 + \beta v^*(2)\} = (2 - \beta)/(1 - \beta) = v^*(2).$$

**Lemma 6.4.** *The refactored  $\sigma$ -value functions are as follows:*

$a$	$g_{\sigma_1}(a)$	$g_{\sigma_2}(a)$	$g_{\sigma_3}(a)$	$g_{\sigma_4}(a)$
0	$\exp[-\gamma/(1 - \beta)]$	$\exp[-\gamma/(1 - \beta)]$	$\exp[-\gamma\beta/(1 - \beta)]$	$\exp[-2\gamma\beta/(1 - \beta^2)]$
1	$\exp[-\gamma(2 - \beta)/(1 - \beta)]$	$\exp[-\gamma/(1 - \beta)]$	$\exp[-\gamma/(1 - \beta)]$	$\exp[-2\gamma/(1 - \beta^2)]$

*Proof.* Regarding  $g_{\sigma_1}$ , by definition,

$$\begin{aligned} g_{\sigma_1}(0) &= \mathbb{E}_{x'|0} \exp \left( -\gamma \left\{ x' - \sigma_1(x') - \frac{\beta}{\gamma} \log g_{\sigma_1}[\sigma_1(x')] \right\} \right) \\ &= \exp \left( -\gamma \left[ 1 - \frac{\beta}{\gamma} \log g_{\sigma_1}(0) \right] \right) = \exp(-\gamma) g_{\sigma_1}(0)^\beta \end{aligned}$$

$$\begin{aligned} \text{and } g_{\sigma_1}(1) &= \mathbb{E}_{x'|1} \exp \left( -\gamma \left\{ x' - \sigma_1(x') - \frac{\beta}{\gamma} \log g_{\sigma_1}[\sigma_1(x')] \right\} \right) \\ &= \exp \left( -\gamma \left[ 2 - \frac{\beta}{\gamma} \log g_{\sigma_1}(0) \right] \right) = \exp(-2\gamma) g_{\sigma_1}(0)^\beta. \end{aligned}$$

So  $g_{\sigma_1}(0)$  and  $g_{\sigma_1}(1)$  are as shown in the table. The remaining proofs are similar.  $\square$

Based on Lemma 6.4, the refactored value function  $g^*$  satisfies

$$g^*(0) = \max_{\sigma_i} g_{\sigma_i}(0) = g_{\sigma_3}(0) = \exp[-\gamma\beta/(1 - \beta)]$$

$$\text{and } g^*(1) = \max_{\sigma_i} g_{\sigma_i}(0) = g_{\sigma_2}(1) = g_{\sigma_3}(1) = \exp[-\gamma/(1 - \beta)].$$



Since  $\hat{g}(0) = W_0 v^*(0) = \exp[-\gamma v^*(1)] = \exp[-\gamma/(1-\beta)] \neq g^*(0)$ , we see that  $g^* \neq \hat{g}$ . Moreover,

$$\begin{aligned} Sg^*(0) &= \exp \left( -\gamma \max \left\{ 1 - \frac{\beta}{\gamma} \log g^*(0), -\frac{\beta}{\gamma} \log g^*(1) \right\} \right) \\ &= \exp[-\gamma(1-\beta+\beta^2)/(1-\beta^2)] \neq g^*(0). \end{aligned}$$

Hence,  $g^* \neq Sg^*$ . The refactored value function is not a fixed point of the refactored Bellman operator.

6.4.2. *Counterexample 2.* We now exhibit a dynamic program and plan factorization under which  $Sg^* = g^*$ ,  $Tv^* \neq v^*$  and  $g^* \neq \hat{g}$ .

The set up is same as Section 6.4.1, except that we let  $\beta > 1$ . In this case, Lemmas 6.3–6.4 still hold and, as in Section 6.4.1, the value function and refactored value function satisfy

$$\begin{aligned} v^*(1) &= 1/(1-\beta), \quad v^*(2) = (2-\beta)/(1-\beta), \\ g^*(0) &= \exp[-\gamma\beta/(1-\beta)] \quad \text{and} \quad g^*(1) = \exp[-\gamma/(1-\beta)]. \end{aligned}$$

We have seen in Section 6.4.1 that  $g^* \neq \hat{g}$ . Since in addition

$$\begin{aligned} Sg^*(0) &= \exp \left( -\gamma \max \left\{ 1 - \frac{\beta}{\gamma} \log g^*(0), -\frac{\beta}{\gamma} \log g^*(1) \right\} \right) \\ &= \exp[-\gamma\beta/(1-\beta)] = g^*(0) \\ \text{and} \quad Sg^*(1) &= \exp \left( -\gamma \max \left\{ 2 - \frac{\beta}{\gamma} \log g^*(0), 1 - \frac{\beta}{\gamma} \log g^*(1) \right\} \right) \\ &= \exp[-\gamma/(1-\beta)] = g^*(1), \end{aligned}$$

we have  $Sg^* = g^*$  as claimed. However, since

$$\begin{aligned} Tv^*(1) &= \max \left\{ 1 - \frac{\beta}{\gamma} \log \mathbb{E}_{x'|0} \exp[-\gamma v^*(x')], -\frac{\beta}{\gamma} \log \mathbb{E}_{x'|1} \exp[-\gamma v^*(x')] \right\} \\ &= \max\{1 + \beta v^*(1), \beta v^*(2)\} = (2\beta - \beta^2)/(1-\beta) \neq v^*(1), \end{aligned}$$

we have  $Tv^* \neq v^*$ , so the value function is not a fixed point of the Bellman operator.

## REFERENCES

BAGGER, J., F. FONTAINE, F. POSTEL-VINAY, AND J.-M. ROBIN (2014): “Tenure, experience, human capital, and wages: A tractable equilibrium search model of wage dynamics,” *American Economic Review*, 104, 1551–96.

- BÄUERLE, N. AND A. JAŚKIEWICZ (2018): “Stochastic optimal growth model with risk sensitive preferences,” *Journal of Economic Theory*, 173, 181–200.
- BELLMAN, R. (1957): *Dynamic programming*, Academic Press.
- BERTSEKAS, D. P. (2012): *Dynamic programming and optimal control*, vol. 2, Athena Scientific, 4th ed.
- (2013): *Abstract dynamic programming*, Athena Scientific.
- BERTSEKAS, D. P. AND H. YU (2012): “Q-learning and enhanced policy iteration in discounted dynamic programming,” *Mathematics of Operations Research*, 37, 66–94.
- BIDDER, R. M. AND M. E. SMITH (2012): “Robust animal spirits,” *Journal of Monetary Economics*, 59, 738–750.
- BLOISE, G. AND Y. VAILAKIS (2018): “Convex dynamic programming with (bounded) recursive utility,” *Journal of Economic Theory*, 173, 118–141.
- DIXIT, A. K. AND R. S. PINDYCK (1994): *Investment Under Uncertainty*, Princeton University Press.
- HANSEN, L. P. AND T. J. SARGENT (2008): *Robustness*, Princeton university press.
- IYENGAR, G. N. (2005): “Robust dynamic programming,” *Mathematics of Operations Research*, 30, 257–280.
- KELLOGG, R. (2014): “The effect of uncertainty on investment: evidence from Texas oil drilling,” *The American Economic Review*, 104, 1698–1734.
- KOCHENDERFER, M. J. (2015): *Decision making under uncertainty: theory and application*, MIT press.
- KRISTENSEN, D., P. MOGENSEN, J. M. MOON, AND B. SCHJERNING (2018): “Solving Dynamic Discrete Choice Models Using Smoothing and Sieve Methods,” Tech. rep., University of Copenhagen.
- LIVSHITS, I., J. MACGEE, AND M. TERTILT (2007): “Consumer bankruptcy: A fresh start,” *American Economic Review*, 97, 402–418.
- LOW, H., C. MEGHIR, AND L. PISTAFERRI (2010): “Wage risk and employment risk over the life cycle,” *American Economic Review*, 100, 1432–67.
- MARINACCI, M. AND L. MONTRUCCHIO (2010): “Unique solutions for stochastic recursive utilities,” *Journal of Economic Theory*, 145, 1776–1804.
- MCCALL, J. J. (1970): “Economics of information and job search,” *The Quarterly Journal of Economics*, 113–126.
- MONAHAN, G. E. (1980): “Optimal stopping in a partially observable Markov process with costly information,” *Operations Research*, 28, 1319–1334.
- MUNOS, R. AND C. SZEPESVÁRI (2008): “Finite-time bounds for fitted value iteration,” *Journal of Machine Learning Research*, 9, 815–857.

- PESKIR, G. AND A. SHIRYAEV (2006): *Optimal stopping and free-boundary problems*, Springer.
- POWELL, W. B. (2007): *Approximate Dynamic Programming: Solving the curses of dimensionality*, vol. 703, John Wiley & Sons.
- PUTERMAN, M. L. AND M. C. SHIN (1982): “Action elimination procedures for modified policy iteration algorithms,” *Operations Research*, 30, 301–318.
- RUST, J. (1987): “Optimal replacement of GMC bus engines: An empirical model of Harold Zurcher,” *Econometrica*, 999–1033.
- (1994): “Structural estimation of Markov decision processes,” *Handbook of Econometrics*, 4, 3081–3143.
- (1996): “Numerical dynamic programming in economics,” *Handbook of Computational Economics*, 1, 619–729.
- RUSZCZYŃSKI, A. (2010): “Risk-averse dynamic programming for Markov decision processes,” *Mathematical Programming*, 125, 235–261.
- SKIENA, S. S. (2008): *The Algorithm Design Manual*, Springer, London.
- TAUCHEN, G. (1986): “Finite state markov-chain approximations to univariate and vector autoregressions,” *Economics letters*, 20, 177–181.