

# DSE 2023 Summer School Lausanne

## Lecture 2: Advances in DP: Theory and Algorithms

John Stachurski

2023

# Topics

## Handling more general dynamic programs

- recursive preferences, quantile preferences, etc.
- recursive preferences + state-dependent discounting, etc.

## Abstraction to simplify proofs

- clarifies optimality proofs
- clarifies relationships between DPs

# Motivation

Consider

$$\max \sum_{t \geq 0} \beta^t u(C_t)$$

subject to

$$W_{t+1} = R(W_t - C_t) \quad \text{and} \quad 0 \leq C_t \leq W_t$$

Standard approach: set up the **Bellman operator**

$$(Tv)(w) = \max_{0 \leq c \leq w} \{u(c) + \beta v(R(w - c))\}$$

# Value function iteration (VFI)

Under some conditions,

1.  $T$  is a contraction mapping
2. the unique fixed point of  $T$  is the value function  $v^*$
3.  $v^*$  can be approximated via  $v^* = \lim_{k \rightarrow \infty} T^k v$  for some  $v$
4. optimal consumption at wealth  $w$  can be found by solving

$$c^* \in \operatorname{argmax}_{0 \leq c \leq w} \{u(c) + \beta v^*(R(w - c))\}$$

# Howard policy iteration

Alternatively, we can use Howard policy iteration (HPI)

A **feasible policy** is a map  $\sigma: \mathbb{R}_+ \rightarrow \mathbb{R}_+$  with

$$0 \leq \sigma(w) \leq w \quad \text{for all } w \in \mathbb{R}_+$$

- given current wealth  $w$ , choose consumption  $c = \sigma(w)$
- $\Sigma :=$  all feasible policies

A feasible policy  $\sigma$  is called  **$v$ -greedy** if

$$\sigma(w) \in \operatorname{argmax}_{0 \leq c \leq w} \{u(c) + \beta v(R(w - c))\}$$

---

**Algorithm 1:** Howard policy iteration

---

input  $\sigma_0 \in \Sigma$ , set  $k \leftarrow 0$  and  $\varepsilon \leftarrow 1$

**while**  $\varepsilon > 0$  **do**

$v_k \leftarrow$  the lifetime value of  $\sigma_k$

$\sigma_{k+1} \leftarrow$  a  $v_k$ -greedy policy

$\varepsilon \leftarrow \mathbb{1}\{\sigma_k \neq \sigma_{k+1}\}$

$k \leftarrow k + 1$

**end**

**return**  $\sigma_k$

---

# Optimistic policy iteration

The lifetime value  $v_\sigma$  of policy  $\sigma$  is the unique  $v$  that solves

$$v(w) = u(\sigma(w)) + \beta v(R(w - \sigma(w)))$$

To compute it we introduce the **policy operator**

$$(T_\sigma v)(w) = u(\sigma(w)) + \beta v(R(w - \sigma(w)))$$

Facts:

1.  $v_\sigma$  is the unique fixed point of  $T_\sigma$
2.  $T_\sigma^k v \rightarrow v_\sigma$  as  $k \rightarrow \infty$  for all reasonable  $v$

---

## Algorithm 2: OPI

---

input  $v_0$ , an initial guess of  $v^*$

input  $\tau$ , a tolerance level for error

input  $m \in \mathbb{N}$ , a step size

$k \leftarrow 0$

$\varepsilon \leftarrow +\infty$

**while**  $\varepsilon > \tau$  **do**

$\sigma_k \leftarrow$  a  $v_k$ -greedy policy

$v_{k+1} \leftarrow T_{\sigma_k}^m v_k$

$\varepsilon \leftarrow \|v_k - v_{k+1}\|_\infty$

$k \leftarrow k + 1$

**end**

**return**  $\sigma_k$

---



## State-dependent discounting

What happens if we introduce state-dependent discounting?

$$(Tv)(w, z) = \max_{0 \leq c \leq w} \left\{ u(c) + \beta(z) \sum_{z'} v(R(w - c), z') Q(z, z') \right\}$$

$T$  is no longer a one-step contraction

However, as discussed in Lecture 1,

- Standard optimality results hold
- OPI, VFI, HPI all converge

## More complications

What happens if we switch to the **expected value function**

$$g(w, z, c) := \sum_{z'} v(R(w - c), z') Q(z, z')$$

with “Bellman operator”

$$(Rg)(w, z, c) =$$

$$\sum_{z'} \max_{0 \leq c' \leq R(w - c)} \{u(c') + \beta(z') g(R(w - c), z', c')\} Q(z, z')$$

Does  $R$  have the same properties as  $T$ ?

What are the equivalent algorithms and do they converge?

And what happens if we introduce **Epstein–Zin preferences**?

$$(Tv)(w, z) =$$

$$\max_{0 \leq c \leq w} \left\{ c^\alpha + \beta(z) \left[ \sum_{z'} v(R(w - c), z')^\gamma Q(z, z') \right]^{\alpha/\gamma} \right\}^{1/\alpha}$$

- Is  $T$  still a contraction?
- Are the previous optimality results still valid?
- Do VFI, OPI, HPI converge?

Or **risk-sensitive preferences**?

$$(Tv)(w, z) =$$

$$\max_{0 \leq c \leq w} \left\{ u(c) + \frac{\beta(z)}{\theta} \ln \left[ \sum_{z'} e^{\theta v(R(w-c), z')} Q(z, z') \right] \right\}$$

- Is  $T$  still a contraction?
- Are the previous optimality results still valid?
- Do VFI, OPI, HPI converge?

What about if we want to handle

- ambiguity?
- expected value VFI in an Epstein–Zin framework?
- expected value VFI + ambiguity + state-dependent discounting?
- integrated value functions in a risk-sensitive framework in continuous time?
- $Q$ -learning?

Is there any unifying theory?

Or are all these problems too diverse?

# Abstraction Level 1: RDPs

1. Construct a DP framework based on an abstraction of Bellman's equation
2. State optimality results in this framework
3. Connect with applications

Builds on work by

- Eric Denardo
- Dimitri Bertsekas
- Takashi Kamihigashi

# Recursive Decision Problems

We begin with a generic version of Bellman's equation:

$$v(x) = \max_{a \in \Gamma(x)} B(x, a, v)$$

- $x \in$  a finite set  $X$  (the **state space**)
- $a \in$  a finite set  $A$  (the **action space**)
- $B(x, a, v)$  = total lifetime rewards
  - contingent on current state-action pair  $(x, a)$  and
  - using  $v$  to evaluate future states

# Definition

A **recursive decision process** (RDP) is a triple  $(\Gamma, V, B)$ , where



- $\Gamma$  is a nonempty correspondence from  $X$  to  $A$

called the **feasible correspondence**

Generates:

- the **feasible state-action pairs**

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

- the set of **feasible policies**

$$\Sigma := \{\sigma \in A^X : \sigma(x) \in \Gamma(x) \text{ for all } x \in X\}$$

- $V$  is a sublattice of  $\mathbb{R}^X$  called the **value space**

- A set of candidates for the value function

(Sublattice  $\iff f \vee g \in V$  and  $f \wedge g \in V$  when  $f, g \in V$ )

- $B$  is a map from  $G \times V$  to  $\mathbb{R}$  called the **value aggregator**

It satisfies

1. **Monotonicity:**

$$v, w \in V \text{ and } v \leq w \implies B(x, a, v) \leq B(x, a, w)$$

for all  $(x, a) \in G$

2. **Consistency:**

$$w(x) := B(x, \sigma(x), v) \text{ is in } V \text{ whenever } \sigma \in \Sigma \text{ and } v \in V$$

**Example.** An arbitrary Markov Decision Process with Bellman equation

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x'} v(x') P(x, a, x') \right\} \quad (1)$$

can be framed as an RDP

Take  $\Gamma$  as given, set  $V = \mathbb{R}^X$ , and

$$B(x, a, v) = r(x, a) + \beta \sum_{x'} v(x') P(x, a, x')$$

- monotonicity and consistency conditions are trivial to check
- recover (1) via  $v(x) = \max_{a \in \Gamma(x)} B(x, a, v)$

**Example.** Consider an **optimal stopping** problem with

$$v(x) = \max \left\{ e(x), c(x) + \beta \sum_{x' \in X} v(x') P(x, x') \right\}$$

Let  $V = \mathbb{R}^X$

If  $\Gamma(x) = \{0, 1\}$  and

$$B(x, a, v) = ae(x) + (1 - a) \left[ c(x) + \beta \sum_{x' \in X} v(x') P(x, x') \right]$$

then  $(\Gamma, V, B)$  is an RDP with the same Bellman equation

**Example.** Consider an MDP with **state-dependent discounting**, so that

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \sum_{x'} v(x') \beta(x, a, x') P(x, a, x') \right\}$$

Let  $V = \mathbb{R}^X$  and

$$B(x, a, v) = r(x, a) + \sum_{x'} v(x') \beta(x, a, x') P(x, a, x')$$

Now  $(\Gamma, V, B)$  is an RDP with the same Bellman equation

**Example.** Consider a modified MDP with **risk-sensitive preferences**, so that

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \frac{1}{\theta} \ln \left( \sum_{x'} \exp(\theta v(x')) P(x, a, x') \right) \right\}$$

for nonzero  $\theta$

With  $V = \mathbb{R}^X$  and

$$B(x, a, v) = r(x, a) + \beta \frac{1}{\theta} \ln \left( \sum_{x'} \exp(\theta v(x')) P(x, a, x') \right)$$

we obtain an RDP with the same Bellman equation

**Example.** Consider a modified MDP with **quantile preferences**, so that

$$v(x) = \max_{a \in \Gamma(x)} \{r(x, a) + \beta(R_\tau^a v)(x)\}$$

where

$$(R_\tau^a v)(x) := \tau\text{-th quantile of } v(X') \text{ when } X' \sim P(x, a, \cdot)$$

With  $V = \mathbb{R}^X$  and

$$B(x, a, v) = r(x, a) + \beta(R_\tau^a v)(x)$$

we obtain an RDP with the same Bellman equation



**Example.** Consider a modified MDP with **Epstein–Zin preferences**, so that

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \left( \sum_{x'} v(x')^\gamma P(x, a, x') \right)^{\alpha/\gamma} \right\}^{1/\alpha}$$

for nonzero  $\alpha, \gamma$

With  $V =$  the strictly positive functions in  $\mathbb{R}^X$  and

$$B(x, a, v) = \left\{ r(x, a) + \beta \left( \sum_{x'} v(x')^\gamma P(x, a, x') \right)^{\alpha/\gamma} \right\}^{1/\alpha}$$

we obtain an RDP with the same Bellman equation

**Example.** Consider a **shortest path problem** on graph  $\mathcal{G} = (\mathbf{X}, E)$

- $c(x, x') = \text{cost of traversing edge } (x, x') \in E$
- the direct successors of  $x$  denoted by

$$\mathcal{O}(x) := \{x' \in \mathbf{X} : (x, x') \in E\}$$

Aim: find the minimum cost path from  $x$  to a specified vertex  $d$

No discounting (so cannot use MDP theory)

The Bellman equation is

$$v(x) = \min_{x' \in \mathcal{O}(x)} \{c(x, x') + v(x')\}$$

Let  $V = \mathbb{R}^X$

Let  $\Gamma(x) = \mathcal{O}(x)$  and

$$B(x, x', v) = c(x, x') + v(x')$$

This is an RDP with the same Bellman equation

# Policies

Consider an arbitrary RDP  $(\Gamma, V, B)$

A **feasible policy** is a

$\sigma \in A^X$  such that  $\sigma(x) \in \Gamma(x)$  for all  $x \in X$

- respond to state  $X_t$  with action  $A_t := \sigma(X_t)$  at **all**  $t \geq 0$
- $\Sigma :=$  the set of all feasible policies

# Policy Operators

Fix  $\sigma \in \Sigma$

The corresponding **policy operator**  $T_\sigma$  is defined at  $v \in V$  by

$$(T_\sigma v)(x) = B(x, \sigma(x), v) \quad (x \in X)$$

**Lemma.**  $T_\sigma$  is an order-preserving self-map on  $V$

Proof: Immediate from monotonicity and consistency

Example. The EZ policy operator is

$$(T_{\sigma} v)(x) = \left\{ r(x, \sigma(x)) + \beta \left( \sum_{x'} v(x')^{\gamma} P(x, \sigma(x), x') \right)^{\alpha/\gamma} \right\}^{1/\alpha}$$

# Optimality

To define optimality for RDPs, we use the natural generalizations...

# Lifetime value

Let  $\mathcal{R} := (\Gamma, V, B)$  be an RDP and let  $\sigma$  be any policy

Suppose  $T_\sigma$  has a unique fixed point in  $V$

We denote this function by  $v_\sigma$  and call it the  **$\sigma$ -value function**

We interpret this function as the lifetime value of following  $\sigma$

We call  $\mathcal{R}$  **well-posed** if  $T_\sigma$  has a unique fixed point in  $V$  for all  $\sigma \in \Sigma$



**Example.** Let  $\mathcal{R}$  be the RDP generated by an MDP

Recall that

$$T_{\sigma} v = r_{\sigma} + \beta P_{\sigma} v$$

This operator has the unique fixed point

$$v_{\sigma} = (I - \beta P_{\sigma})^{-1} r_{\sigma}$$

- Hence  $\mathcal{R}$  is well-posed
- $v_{\sigma}(x) = \mathbb{E}_x \sum_{t \geq 0} \beta^t r(X_t, \sigma(X_t)) = \text{lifetime value}$

Example. For the Epstein–Zin RDP,

$$(T_{\sigma}v)(x) = \left\{ r(x, \sigma(x)) + \beta \left[ \sum_{x' \in X} v(x')^{\gamma} P(x, \sigma(x), x') \right]^{\alpha/\gamma} \right\}^{1/\alpha}$$

and

$V :=$  the strictly positive functions in  $\mathbb{R}^X$

- Is this RDP well-posed?

# Greedy Policies

Fix  $v \in \mathbb{R}^X$

A policy  $\sigma$  is called  **$v$ -greedy** if

$$\sigma(x) \in \operatorname{argmax}_{a \in \Gamma(x)} B(x, a, v)$$

for all  $x \in X$

Note: at least one  $v$ -greedy policy exists in  $\Sigma$

# The Bellman Operator

The **Bellman operator** is the self-map on  $\mathbb{R}^X$  defined by

$$(Tv)(x) = \max_{a \in \Gamma(x)} B(x, a, v)$$

Key idea:

$$Tv = v \iff v \text{ satisfies the Bellman equation}$$

# Optimality

Let  $\mathcal{R}$  be a well-posed RDP

The **value function** is defined by  $v^* = \bigvee v_\sigma$

More explicitly,

$$v^*(x) := \max_{\sigma \in \Sigma} v_\sigma(x) \quad (x \in X)$$

= max lifetime value from state  $x$

A policy  $\sigma \in \Sigma$  is called **optimal** if

$$v_\sigma = v^*$$

## Howard policy iteration for RDPs

---

---

```
input  $\sigma_0 \in \Sigma$ , an initial guess of  $\sigma^*$ 
 $k \leftarrow 0$ 
 $\varepsilon \leftarrow 1$ 
while  $\varepsilon > 0$  do
     $v_k \leftarrow$  the unique fixed point of  $T_{\sigma_k}$ 
     $\sigma_{k+1} \leftarrow$  a  $v_k$  greedy policy
     $\varepsilon \leftarrow \mathbb{1}\{\sigma_k \neq \sigma_{k+1}\}$ 
     $k \leftarrow k + 1$ 
end
return  $\sigma_k$ 
```

---

Let  $\mathcal{R}$  be an RDP

Key question:

What assumptions to we need for optimality?

Obviously  $\mathcal{R}$  must be well-posed

- each  $T_\sigma$  has a unique fixed point in  $V$

This is the minimum requirement

What else?

# Stability

Let  $\mathcal{R}$  be an RDP

We call  $\mathcal{R}$  **globally stable** if, for all  $\sigma \in \Sigma$ , the operator  $T_\sigma$  is globally stable on  $V$

Meaning

1.  $T_\sigma$  has a unique fixed point in  $V$  and
2.  $\lim_{k \rightarrow \infty} T_\sigma^k v = v_\sigma$  for all  $v \in V$



Let  $\mathcal{R}$  be a well-posed RDP with value function  $v^*$

**Theorem.** If  $\mathcal{R}$  is stable, then

1.  $v^*$  is the unique solution to the Bellman equation in  $\mathbb{R}^X$
2. A feasible policy is optimal if and only if it is  $v^*$ -greedy
3. At least one optimal policy exists
4. Howard policy iteration returns an optimal policy in finitely many steps

# Types of RDPs

The optimality properties require stability

We can check this directly

We can also

1. identify classes of RDPs that are stable
2. show that a given application belongs to one of these classes

Let's discuss the classification approach

Below  $\mathcal{R} = (\Gamma, V, B)$  is a fixed RDP

# Contracting RDPs

We call  $\mathcal{R}$  **contracting** if  $\exists \beta < 1$  such that

$$|B(x, a, v) - B(x, a, w)| \leq \beta \|v - w\|_{\infty}$$

for all  $(x, a) \in G$  and  $v, w \in V$

**Thm.** If  $\mathcal{R}$  is contracting and  $V$  is closed, then  $\mathcal{R}$  is stable

Proof: Easy to show that each  $T_{\sigma}$  is a contraction on  $V$

(Main idea dates back to Denardo 1967)

## Eventually Contracting RDPs

We call  $\mathcal{R}$  **eventually contracting** if there is an  $L \geq 0$  such that  $\rho(L) < 1$  and

$$|B(x, a, v) - B(x, a, w)| \leq \sum_{x'} |v(x') - w(x')| L(x, x')$$

for all  $(x, a) \in G$  and  $v, w \in V$

**Thm.** If  $\mathcal{R}$  is eventually contracting and  $V$  is closed, then  $\mathcal{R}$  is stable

Proof: See the book

# Concave RDPs

We call  $\mathcal{R}$  **concave** if

1.  $V = [v_1, v_2]$
2.  $B(x, a, v_1) > v_1(x)$  for all  $(x, a) \in G$  and
3.  $v \mapsto B(x, a, v)$  is concave for all  $(x, a) \in G$

**Thm.** If  $\mathcal{R}$  is concave, then  $\mathcal{R}$  is stable

Proof: See the book

# Application: job search with quantile preferences

Set up:

- wage offer process  $(W_t)_{t \geq 0}$  is  $P$ -Markov on finite set  $W$
- discount factor  $\beta \in (0, 1)$

The Bellman equation is

$$v(w) = \max \left\{ \frac{w}{1 - \beta}, c + \beta(R_\tau v)(w) \right\}$$

Here

$$(R_\tau v)(w) := \tau\text{-th quantile of } v(W') \text{ when } W' \sim P(w, \cdot)$$

This problem studied in

- de Castro and Galvao (2019)
- de Castro, Galvao and Nunes (2022)
- de Castro and Galvao (2022)

We can embed the into the RDP framework by taking

- $\Gamma(w) = \{0, 1\}$
- $V = \mathbb{R}_+^W$
- $B$  given by

$$B(w, a, v) = a \frac{w}{1 - \beta} + (1 - a)[c + \beta(R_\tau v)(w)]$$

Easy to check that  $\mathcal{R} := (\Gamma, V, B)$  is an RDP with Bellman equation

$$v(w) = \max \left\{ \frac{w}{1 - \beta}, c + \beta(R_\tau v)(w) \right\}$$

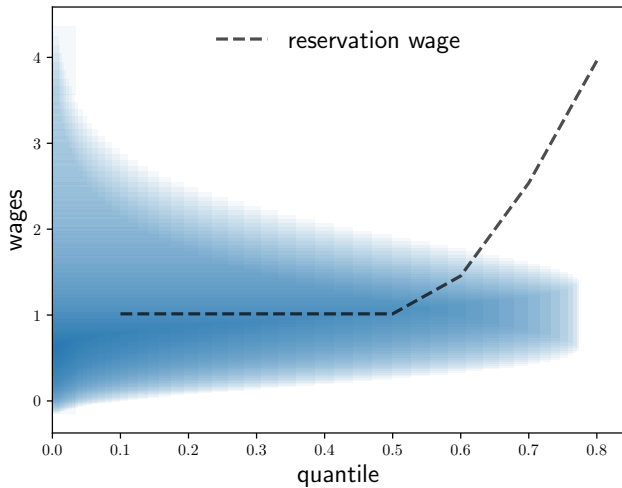


**Proposition.**  $\mathcal{R}$  is a contracting RDP

- Proof: See the text

Since  $V$  is closed,  $\mathcal{R}$  is stable

Hence all optimality properties apply



## Abstraction Level 2: ADPs

We define an **abstract dynamic program (ADP)** to be a pair

$$\mathcal{A} = (V, \{T_\sigma\}_{\sigma \in \Sigma}), \quad \text{where}$$

1.  $V = (V, \preceq)$  is a partially ordered set and
2.  $\{T_\sigma\}_{\sigma \in \Sigma}$  is a family of self-maps on  $V$

Below,

- elements of  $\Sigma$  will be referred to as **policies**
- elements of  $\{T_\sigma\}$  are called **policy operators**

If  $T_\sigma$  has a unique fixed point, then we

- denote it  $v_\sigma$  and call it the  **$\sigma$ -value function**
- understand  $v_\sigma$  as representing lifetime value of  $\sigma$

Interpretation:

- $V$  is a set of candidate value functions
- $\Sigma$  is a set of feasible policies
- the lifetime value of  $\sigma \in \Sigma$  is  $v_\sigma$
- we seek a greatest element in  $\{v_\sigma\}_{\sigma \in \Sigma}$

**Example.** Consider an RDP  $(\Gamma, V, B)$

Let  $\Sigma$  be the set of feasible policies

Recall that, for each  $\sigma \in \Sigma$ , the policy operator  $T_\sigma$  is defined at  $v \in V$  by

$$(T_\sigma v)(x) = B(x, \sigma(x), v)$$

The pair  $(V, \{T_\sigma\})$  is an ADP

# Benefits of ADP theory

- More abstraction means easier proofs
- Removing structure makes it easier to see connections
- Can handle a more diverse range of problems

Given  $v \in V$ , a policy  $\sigma$  in  $\Sigma$  is called  **$v$ -greedy** if

$$T_\sigma v \succeq T_\tau v \quad \text{for all } \tau \in \Sigma$$

**Example.** In the MDP example we have

$$(T_\sigma v)(x) = r(x, \sigma(x)) + \beta \sum_{x'} v(x') P(x, \sigma(x), x')$$

so  $\sigma$  is  $v$ -greedy iff

$$\sigma(x) \in \operatorname{argmax}_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x'} v(x') P(x, a, x') \right\} \quad \text{for all } x \in X$$

# Bellman equation

Fix an ADP  $\mathcal{A} = (V, \{T_\sigma\})$

We define the **Bellman operator** via

$$Tv := \bigvee_{\sigma} T_{\sigma} v$$

(if it exists)

Equivalently,

$$Tv = T_{\sigma} v \text{ when } \sigma \text{ is } v\text{-greedy}$$

We say that  $v \in V$  satisfies the **Bellman equation** if  $Tv = v$



Example. For the MDP,

$(Tv)(x) = (T_\sigma v)(x)$  when  $\sigma$  is  $v$ -greedy

$$= \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x'} v(x') P(x, a, x') \right\}$$

Hence the ADP Bellman equation is

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \beta \sum_{x'} v(x') P(x, a, x') \right\}$$

And this is the same as the MDP Bellman equation

# Properties

We say that  $\mathcal{A} = (V, \{T_\sigma\})$  is

- **well-posed** if  $T_\sigma$  has one fixed point in  $V$  for each  $\sigma \in \Sigma$
- **order stable** if  $(V, T_\sigma)$  is order stable for each  $\sigma \in \Sigma$
- **max-stable** if  $\mathcal{A}$  is order stable, each  $v \in V$  has at least one greedy policy, and  $T$  has at least one fixed point in  $V$

Note: order stability is a regularity property — see Ch 9

Let  $\mathcal{A}$  be a well-posed ADP

A policy  $\sigma \in \Sigma$  is called **optimal** for  $\mathcal{A}$  if

$$v_\tau \preceq v_\sigma \text{ for all } \tau \in \Sigma$$

We set  $v^* := \bigvee_\sigma v_\sigma$  and call  $v^*$  the **value function**

We define a self-map  $H$  on  $V$  via

$$H v = v_\sigma \quad \text{where } \sigma \text{ is } v\text{-greedy}$$

Iterating with  $H$  is an abstract version of HPI

# Max-Optimality

**Theorem.** If  $\mathcal{A}$  is max-stable, then

1.  $v^*$  exists in  $V$
2.  $v^*$  is the unique solution to the Bellman equation in  $V$
3. a policy is optimal if and only if it is  $v^*$ -greedy
4. at least one optimal policy exists

If, in addition,  $\Sigma$  is finite, then  $\text{HPI} \rightarrow v^*$  in finitely many steps

## Subordinate ADPs

Let  $\mathcal{A} := (V, \{T_\sigma\})$  and  $\hat{\mathcal{A}} := (\hat{V}, \{\hat{T}_\sigma\})$  be ADPs

We say that  $\hat{\mathcal{A}}$  is **subordinate** to  $\mathcal{A}$  if  $\exists$

1. an order-preserving map  $F$  from  $V$  onto  $\hat{V}$  and
2. order-preserving maps  $\{G_\sigma\}_{\sigma \in \Sigma}$  from  $\hat{V}$  to  $V$

such that

$$T_\sigma = G_\sigma \circ F \quad \text{and} \quad \hat{T}_\sigma = F \circ G_\sigma \quad \text{for all } \sigma \in \Sigma$$

Let  $G = \bigvee_\sigma G_\sigma$

**Theorem.** If

1.  $\mathcal{A}$  is max-stable and
2.  $\hat{\mathcal{A}}$  is subordinate to  $\mathcal{A}$ ,

then  $\hat{\mathcal{A}}$  is also max-stable and the Bellman operators are related by

$$T = G \circ F \quad \text{and} \quad \hat{T} = F \circ G$$

while the value functions are related by

$$v^* = G \hat{v}^* \quad \text{and} \quad \hat{v}^* = F v^*$$

Moreover,

1. if  $\sigma$  is optimal for  $\mathcal{A}$ , then  $\sigma$  is optimal for  $\hat{\mathcal{A}}$ , and
2. if  $G_\sigma \hat{v}^* = G \hat{v}^*$ , then  $\sigma$  is optimal for  $\mathcal{A}$

# Application

Consider an Epstein–Zin dynamic program with Bellman equation

$$v(w, e) = \max_{0 \leq s \leq w} \left\{ r(w, s, e)^\alpha + \beta \left( \sum_{e'} v(s, e')^\gamma \varphi(e') \right)^{\alpha/\gamma} \right\}^{1/\alpha}$$

Here

- $w$  is current wealth (discretized)
- $s$  is savings (discretized)
- $e$  is an IID endowment shock with range  $E$
- $\beta$  is a constant in  $(0, 1)$  and  $r$  is a reward function

The policy operator corresponding to  $\sigma \in \Sigma$  is

$$(T_\sigma v)(w, e) = \left\{ r(w, \sigma(w), e)^\alpha + \beta \left( \sum_{e'} v(\sigma(w), e')^\gamma \varphi(e') \right)^{\alpha/\gamma} \right\}^{1/\alpha}$$

**Proposition.** If

- $X := W \times E$  and
- $V := (0, \infty)^X$ ,

then  $\mathcal{A} = (V, \{T_\sigma\})$  is a max-stable ADP

(Details in Ch 9)



Next consider the operator

$$(B_\sigma h)(w) = \left\{ \sum_e \{r(w, \sigma(w), e)^\alpha + \beta h(\sigma(w))^\alpha\}^{\gamma/\alpha} \varphi(e) \right\}^{1/\gamma},$$

where  $h$  is an element of  $(0, \infty)^W$

Define  $F$  at  $v \in V$  by

$$(Fv)(w) = \left\{ \sum_e v(w, e)^\gamma \varphi(e) \right\}^{1/\gamma} \quad (w \in W)$$

Then  $\mathcal{B} = (F(V), \{B_\sigma\})$  is also an ADP

Moreover,  $\mathcal{B}$  is subordinate to  $\mathcal{A}$

To see, this, define  $G_\sigma$  by

$$(G_\sigma h)(w, e) = \{r(w, \sigma(w), e)^\alpha + \beta h(\sigma(w))^\alpha\}^{1/\alpha}$$

Then

- $F$  and  $G_\sigma$  are order-preserving
- $T_\sigma$  is equal to  $G_\sigma \circ F$  and
- $B_\sigma$  is equal to  $F \circ G_\sigma$

---

**Algorithm 3:** Solving  $\mathcal{A}$  via  $\mathcal{B}$ 

---

input  $\sigma_0 \in \Sigma$ , set  $k \leftarrow 0$  and  $\varepsilon \leftarrow 1$

**while**  $\varepsilon > 0$  **do**

$h_k \leftarrow$  the fixed point of  $B_{\sigma_k}$

$\sigma_{k+1} \leftarrow$  an  $h_k$ -greedy policy, satisfying

$$\sigma_{k+1}(w) \in \operatorname{argmax}_{0 \leq s \leq w} \left\{ \sum_e \{r(w, s, e)^\alpha + \beta h(s)^\alpha\}^{\gamma/\alpha} \varphi(e) \right\}^{1/\gamma}$$

$\varepsilon \leftarrow \mathbb{1}\{\sigma_k \neq \sigma_{k+1}\}$  and  $k \leftarrow k + 1$

**end**

Compute  $\sigma$  to satisfy

$$\sigma(w, e) \in \operatorname{argmax}_{0 \leq s \leq w} \{r(w, s, e)^\alpha + \beta h_k(s)^\alpha\}^{1/\alpha}$$

**return**  $\sigma$

---

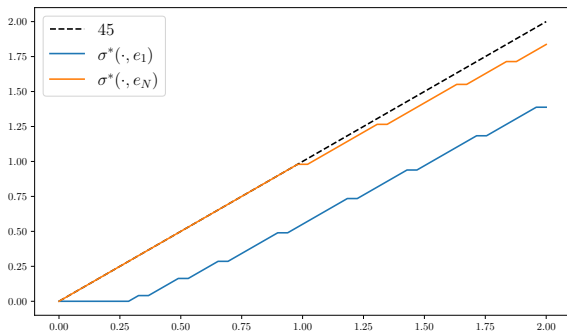


Figure: Optimal savings policy with Epstein–Zin preference

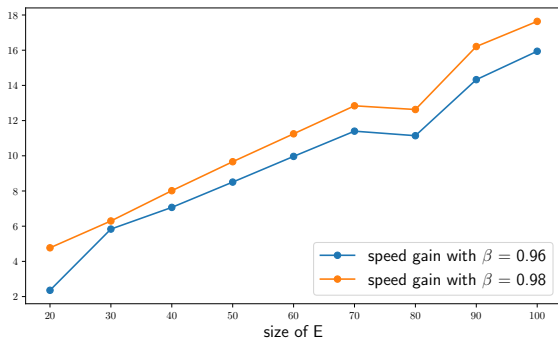


Figure: Speed gain from replacing  $\mathcal{A}$  with subordinate model  $\mathcal{B}$

For details of computations see

[https://github.com/jstac/adps\\_public](https://github.com/jstac/adps_public)