

Advances in Dynamic Programming

Theory and Applications

Thomas J. Sargent and John Stachurski

2023

Plan

1. Review traditional dynamic programs (MDPs)
2. Generalize to recursive decision problems (RDPs)
3. Discuss RDP optimality
4. Identify types of RDPs
5. Consider some applications

RDPs are a generalization of MDPs that can handle

- recursive preferences
- robust control
- ambiguity
- state-dependent discounting
- adversarial agents
- negative discount rates
- jump processes (continuous time)
- etc., etc.

Aims

Provide conditions under which

- the value function satisfies Bellman's equation
- Bellman's principle of optimality holds
- at least one optimal policy exists
- standard algorithms converge
 - value function iteration (VFI)
 - optimistic policy iteration (OPI)
 - Howard policy iteration (HPI)

The RDP framework builds on work by

- Eric Denardo
- Dimitri Bertsekas
- Takashi Kamihigashi
- etc.

Further references

- [Dynamic Programming \(Vol 1\)](#)
- [Completely Abstract Dynamic Programming](#)

Terminology

Let T be a self-map on metric space $U = (U, \rho)$

We call T **globally stable** on U when

1. T has a unique fixed point \bar{v} in V and
2. $T^k v \rightarrow \bar{v}$ as $k \rightarrow \infty$ for all $v \in V$

Let T be a self-map on partially ordered space $U = (U, \leq)$

We call T **order preserving** on U when

1. $u, v \in U$ and $u \leq v$ implies $Tu \leq Tv$

Fixed point theory: quick reminders

Our old friend **Banach**

Theorem. If

1. (U, d) is a complete metric space
2. T is a contraction of modulus λ on U

then T has a unique fixed point u^* in U and

$$d(T^k u, u^*) \leq \lambda^k d(u, u^*) \quad \text{for all } k \in \mathbb{N} \text{ and } u \in U$$

In particular, T is globally stable on U

An “eventual contraction” result:

- X is finite and T is a self-map on $V \subset \mathbb{R}^X$
- V is closed

Thm. T is globally stable on $V \subset \mathbb{R}^X$ whenever there exists a positive linear operator L on \mathbb{R}^X such that

1. $\rho(L) < 1$ and
2. for all $u, v \in V$, we have

$$|Tu - Tv| \leq L|u - v|$$

(For a proof see Ch. 6)

Du's Theorem

Let

- X be a finite set
- $I := [v_1, v_2]$ be a nonempty order interval in (\mathbb{R}^X, \leq)
- T be an order-preserving self-map on I

Thm. T is globally stable on I if either

1. T is concave and $Tv_1 \gg v_1$ or
2. T is convex and $Tv_2 \ll v_2$

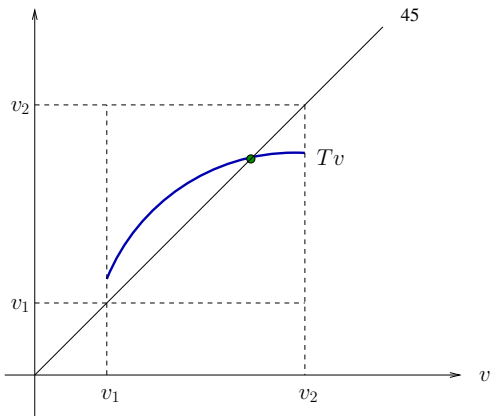


Figure: Concave case (one-dimensional)

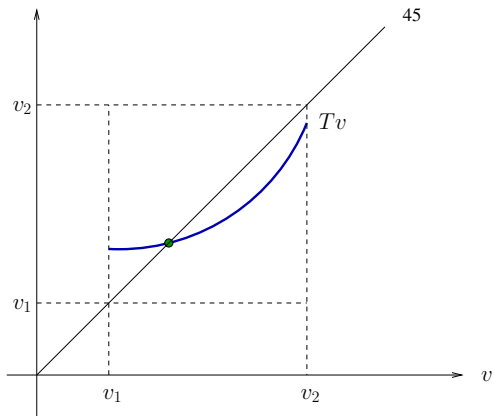


Figure: Convex case (one-dimensional)

Review

Let's quickly review Markov decision processes (MDPs)

- the most standard framework for dynamic programming
- includes many standard economic applications
- useful for intuition
- but limited in some ways

Aims:

- Recall basic concepts
- Lay groundwork for discussing RDPs

States and Actions

We take as given

1. a finite set X called the **state space** and
2. a finite set A called the **action space**

Actions are restricted by a **feasible correspondence** Γ from X to A

- $\Gamma(x)$ = actions available in state x (nonempty)

Given Γ , we define the **feasible state-action pairs**

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

Dynamics

A **stochastic kernel** from G to X is a map

$$P: G \times X \rightarrow [0, 1]$$

satisfying

$$\sum_{x'} P(x, a, x') = 1 \quad \text{for all } (x, a) \text{ in } G$$

- next period state x' is drawn from $P(x, a, \cdot)$

Rewards

Flow reward $r(x, a)$ is received at $(x, a) \in G$

Lifetime rewards are

$$\mathbb{E} \sum_{t \geq 0} \beta^t r(X_t, A_t)$$

where $\beta \in (0, 1)$

- β is called the **discount factor**
- r is called the **reward function**

The **Bellman equation** is

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

Formal definition

Let's summarize the primitives:

Given X and A , an **MDP** is a tuple (Γ, β, r, P) where

1. Γ is a nonempty correspondence from $X \rightarrow A$
2. β is a constant in $(0, 1)$
3. r is a function from G to \mathbb{R}
4. P is a stochastic kernel from G to X

Policies

A **feasible policy** is a $\sigma \in A^X$ such that

$$\sigma(x) \in \Gamma(x) \text{ for all } x \in X$$

Choose $\sigma \iff$

respond to state X_t with action $\sigma(X_t)$ at all $t \geq 0$

Notation:

$\Sigma :=$ the set of all feasible policies

Closed loop dynamics

Choosing $\sigma \in \Sigma \implies$

1. flow rewards at x are

$$r_\sigma(x) := r(x, \sigma(x))$$

2. the state process $(X_t)_{t \geq 0}$ follows

$$P_\sigma(x, x') := P(x, \sigma(x), x')$$

(We say that $(X_t)_{t \geq 0}$ is **P_σ -Markov**)

If

$(X_t)_{t \geq 0}$ is P_σ -Markov with $X_0 = x$

then the **lifetime value of σ** starting from x is

$$v_\sigma(x) := \mathbb{E}_x \sum_{t \geq 0} \beta^t r(X_t, \sigma(X_t))$$

$$= \mathbb{E}_x \sum_{t \geq 0} \beta^t r_\sigma(X_t)$$

$$= \sum_{t \geq 0} \beta^t \mathbb{E}_x r_\sigma(X_t)$$

$$= \sum_{t \geq 0} \beta^t (P_\sigma^t r_\sigma)(x)$$

In vector notation,

$$v_{\sigma} = \sum_{t \geq 0} (\beta P_{\sigma})^t r_{\sigma}$$

Since $\beta < 1$,

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

Policy Operators

Given $\sigma \in \Sigma$, the **policy operator** T_σ is defined by

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

In vector notation,

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

Lemma. T_σ is order-preserving on \mathbb{R}^X

Proof: $v \leq w \implies P_\sigma v \leq P_\sigma w \implies T_\sigma v \leq T_\sigma w$

Lemma. T_σ is a contraction of modulus β on \mathbb{R}^X

Proof: Let $|v| := v \vee (-v)$ and fix v, w in \mathbb{R}^X

We have

$$\begin{aligned}|T_\sigma v - T_\sigma w| &= \beta |P_\sigma v - P_\sigma w| \\ &\leq \beta P_\sigma |v - w| \\ &\leq \beta P_\sigma \|v - w\|_\infty \mathbb{1} \\ &= \beta \|v - w\|_\infty \mathbb{1}\end{aligned}$$

Finally, $|a| \leq |b|$ implies $\|a\|_\infty \leq \|b\|_\infty$

Lemma. T_σ is a contraction of modulus β on \mathbb{R}^X

Proof: Let $|v| := v \vee (-v)$ and fix v, w in \mathbb{R}^X

We have

$$\begin{aligned} |T_\sigma v - T_\sigma w| &= \beta |P_\sigma v - P_\sigma w| \\ &\leq \beta P_\sigma |v - w| \\ &\leq \beta P_\sigma \|v - w\|_\infty \mathbb{1} \\ &= \beta \|v - w\|_\infty \mathbb{1} \end{aligned}$$

Finally, $|a| \leq |b|$ implies $\|a\|_\infty \leq \|b\|_\infty$

Key idea: Lifetime value of $\sigma \iff$ fixed point of T_σ

Indeed,

$$v = T_\sigma v \iff v = r_\sigma + \beta P_\sigma v$$

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

$$\iff v = v_\sigma$$

Greedy Policies

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

A policy σ is called **v -greedy** if

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

for all $x \in X$

Note: at least one v -greedy policy exists in Σ

The Bellman Operator

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

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

- $Tv = v \iff v$ satisfies the Bellman equation

Note

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

Equivalently, $Tv = \bigvee_{\sigma} T_{\sigma} v$

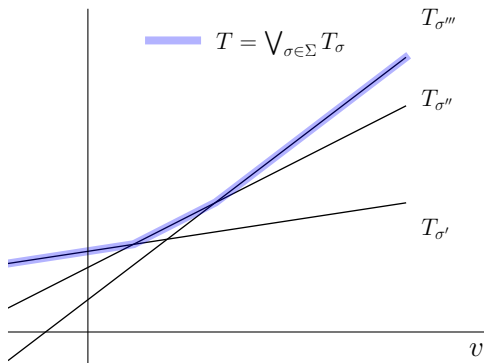


Figure: Visualization in one dimension

Theorem. T is globally stable on \mathbb{R}^X

Proof: Easy to check that T is a contraction of modulus β

Optimality

The **value function** $v^* \in \mathbb{R}^X$ is defined by

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

= max lifetime value from state x

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

$$v_{\sigma} = v^*$$

Howard policy iteration (HPI)

input $\sigma_0 \in \Sigma$

$k \leftarrow 0$

repeat

$v_k \leftarrow (I - \beta P_{\sigma_k})^{-1} r_{\sigma_k}$

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

$k \leftarrow k + 1$

until $\sigma_k = \sigma_{k-1}$

return σ_k

Theorem. For any MDP with value function v^* ,

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. HPI returns an exact optimal policy in finitely many steps

Remark: Point (2) is called **Bellman's principle of optimality**

Where to now?

Researchers are pushing past the boundaries of MDPs

Problems that do not fit the MDP framework include

1. models with nonlinear recursive preferences
2. models with stochastic discounting
3. recursive equilibria in economic geography, production, etc.
4. problems with ambiguity, adversarial agents
5. various combinations of the above, etc.

Where to now?

Researchers are pushing past the boundaries of MDPs

Problems that do not fit the MDP framework include

1. models with nonlinear recursive preferences
2. models with stochastic discounting
3. recursive equilibria in economic geography, production, etc.
4. problems with ambiguity, adversarial agents
5. various combinations of the above, etc.

Our plan

1. Construct an abstract DP framework that includes MDPs as a special case
2. State optimality results in this framework
3. Provide sufficient conditions for several important cases
4. Connect with applications

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**)
- v is a candidate value function
- $B(x, a, v)$ = total lifetime rewards given x, a, v

More formally...

A **recursive decision process** (RDP) is a triple (Γ, V, B) , where

1. Γ is a nonempty correspondence from X to A

- called the **feasible correspondence**
- generates the **feasible state-action pairs**

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

2. V is a nonempty subset of \mathbb{R}^X

- called the **value space**
- A set of candidates for the value function

3. B is a map from $G \times V$ to \mathbb{R} satisfying

(a) **monotonicity**:

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

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

(b) **consistency**:

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

- B is called the **value aggregator**

Example. Consider an MDP with Bellman equation

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

Proof: Take Γ as is, $V := \mathbb{R}^X$ and

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

Now (Γ, V, B) is an RDP

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

Example. Optimal stopping

$$v(x) = \max \left\{ e(x), c(x) + \beta \sum_{x'} v(x') P(x, x') \right\} \quad (2)$$

Set $V := \mathbb{R}^X$, $\Gamma(x) := \{0, 1\}$ and

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

Then (Γ, V, B) is an RDP — checking conditions is trivial

- Setting $v(x) = \max_{a \in \Gamma(x)} B(x, a, v)$ recovers (2)

Example. State-dependent discounting

$$v(x) = \max_{a \in \Gamma(x)} \left\{ r(x, a) + \sum_{x'} v(x') \beta(x, a, x') P(x, a, x') \right\} \quad (3)$$

Take Γ as is, $V := \mathbb{R}^X$ and

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

- (Γ, V, B) is an RDP
- Setting $v(x) = \max_{a \in \Gamma(x)} B(x, a, v)$ recovers (2)

Example. Risk-sensitive preferences

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

for nonzero θ

Take Γ as is, $V := \mathbb{R}^X$ and

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

- an RDP with Bellman equation $v(x) = \max_{a \in \Gamma(x)} B(x, a, v)$

Example. Quantile preferences

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

where

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

Take Γ as is, $V := \mathbb{R}^X$ and

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

- an RDP with Bellman equation $v(x) = \max_{a \in \Gamma(x)} B(x, a, v)$

Example. Epstein–Zin preferences

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

for nonzero α, γ and $r \geq 0$

Take Γ as is, $V := (0, \infty)^X$ and

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

- an RDP with Bellman equation $v(x) = \max_{a \in \Gamma(x)} B(x, a, v)$

Example. **Shortest path problem** on digraph $\mathcal{G} = (\mathbf{X}, E)$

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

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

Aim:

find minimum cost path from x to destination $d \in \mathbf{X}$

No discounting — not an MDP

The Bellman equation is

$$v(x) = \min_{x' \in \mathcal{O}(x)} \{c(x, x') + v(x')\} \quad (4)$$

Set $V := \mathbb{R}^X$, $\Gamma(x) := \mathcal{O}(x)$ and

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

Then (Γ, V, B) is an RDP and

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

recovers (4) — this is minimization, which we treat in Ch. 9

RDP setting can also handle

- ambiguity (smooth ambiguity model)
- adversarial agents
- negative discount rates
- recursive equilibria in economic geography, production, etc.
- various combinations of the above, etc.
- jump processes (continuous time)

See Ch.s 8–10 and discussion below

RDP theory

So far we have just defined RDPs — little structure

Next steps

1. define lifetime values
2. define optimality
3. provide conditions for optimality
4. discuss algorithms

Policies

Fix arbitrary RDP $\mathcal{R} = (\Gamma, V, B)$

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 with action $a := \sigma(x)$ at all $t \geq 0$
- $\Sigma :=$ the set of all feasible policies

Policy Operators

Fix $\sigma \in \Sigma$

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

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

Lemma. T_σ is an order-preserving self-map on V

Proof: Immediate from monotonicity and consistency

Example. The EZ policy operator is

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

Example. The state-dependent discounting policy operator is

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

Lifetime value

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

If T_σ has a unique fixed point in V we denote it by v_σ

- Interpretation: v_σ is the lifetime value of following σ

We call \mathcal{R} **well-posed** if

T_σ has a unique fixed point in V for all $\sigma \in \Sigma$

- A minimal condition for discussing optimality

Why is the “lifetime value” interpretation valid?

Example. Let \mathcal{R} be the RDP generated by MDP (Γ, β, r, P)

In this case

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

Note \mathcal{R} is well-posed because T_σ has unique fixed point

$$v_\sigma = \sum_{t \geq 0} \beta^t P^t r_\sigma = (I - \beta P_\sigma)^{-1} r_\sigma$$

Moreover,

$$v_\sigma(x) = \mathbb{E}_x \sum_{t \geq 0} \beta^t r(X_t, \sigma(X_t)) = \text{lifetime value under } \sigma$$

Example. Consider the Epstein–Zin RDP

A fixed point of T_σ obeys

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

A fixed point of this equation is how we define lifetime value from each state under EZ preferences

- Is this RDP well-posed?

Greedy Policies

Given $v \in \mathbb{R}^X$, a policy σ is called **v -greedy** if

$$\sigma(x) \in \operatorname{argmax}_{a \in \Gamma(x)} B(x, a, v) \quad \text{for all } x \in X$$

- Note: at least one v -greedy policy exists in Σ

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

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

- Note: $Tv = v \iff v$ satisfies the Bellman equation

Note that

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

Proof: Fix $v \in V$

By definition, σ is v -greedy if and only if

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

This is equivalent to

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

Note that

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

Proof: Fix $v \in V$

By definition, σ is v -greedy if and only if

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

This is equivalent to

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

Optimality

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

The **value function** v^* is defined by

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

= max lifetime value from state x

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

$$v_{\sigma} = v^*$$

HPI for well-posed RDPs

input $\sigma_0 \in \Sigma$

$k \leftarrow 0$

repeat

$v_k \leftarrow$ the unique fixed point of T_{σ_k}

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

$k \leftarrow k + 1$

until $\sigma_k = \sigma_{k-1}$

return σ_k

Let \mathcal{R} be an RDP

Key question:

What assumptions to we need for optimality?

Obviously \mathcal{R} must be well-posed

- each σ must have a uniquely defined lifetime value

This is the minimum requirement

What else do we need?

We call \mathcal{R} **globally stable** if

T_σ is globally stable on V for all $\sigma \in \Sigma$

Let \mathcal{R} be an RDP

Key question:

What assumptions to we need for optimality?

Obviously \mathcal{R} must be well-posed

- each σ must have a uniquely defined lifetime value

This is the minimum requirement

What else do we need?

We call \mathcal{R} **globally stable** if

T_σ is globally stable on V for all $\sigma \in \Sigma$

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

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

1. v^* is in V
2. v^* is the unique solution to the Bellman equation in V
3. A feasible policy is optimal if and only if it is v^* -greedy
4. At least one optimal policy exists
5. HPI returns an exact optimal policy in finitely many steps

Proof: See Ch. 9

Types of RDPs

The key condition above is global stability

We can check this directly — show each T_σ is globally stable

We can also

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

Let's discuss the second approach

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

Contracting RDPs

We call \mathcal{R} **contracting** if

$$\exists \beta < 1 \text{ such that } |B(x, a, v) - B(x, a, w)| \leq \beta \|v - w\|_\infty$$

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 globally stable

- Hence all optimality results on slide 58 hold

Proof: Easy to show each T_σ is a mod- β contraction on $(V, \|\cdot\|_\infty)$

(Main idea dates back to Denardo 1967)

Eventually Contracting RDPs

We call \mathcal{R} **eventually contracting** if \exists an $L \geq 0$ s.t.

1. $\rho(L) < 1$ and
2. for all $(x, a) \in G$ and $v, w \in V$,

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

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

- Hence all optimality results on slide 58 hold

Proof: Let \mathcal{R} be eventually contracting

Fixing $\sigma \in \Sigma$ and $v, w \in V$,

$$\begin{aligned} |(T_\sigma v)(x) - (T_\sigma w)(x)| &= |B(x, \sigma(x), v) - B(x, \sigma(x), w)| \\ &\leq \sum_{x'} |v(x') - w(x')| L(x, x') \end{aligned}$$

In other words,

$$|T_\sigma v - T_\sigma w| \leq L|v - w|$$

The claim now follows from the result on slide 8

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 globally stable

- Hence all optimality results on slide 58 hold

Proof: Let \mathcal{R} be concave and fix $\sigma \in \Sigma$

Recall that T_σ is an order-preserving self-map on V

Also, $v \mapsto B(x, \sigma(x), v) = (T_\sigma v)(x)$ is concave for all $x \in X$

Hence T_σ is a concave operator

Also, $T_\sigma v_2 \leq v_2$ because T_σ is a self-map on $V = [v_1, v_2]$

Finally, for any $x \in X$,

$$(T_\sigma v_1)(x) = B(x, \sigma(x), v_1) > v_1(x)$$

Now apply Du's theorem on slide 9

Applications

We have just listed three classes of globally stable RDPS

1. contracting RDPS
2. eventually contracting RDPS
3. concave RDPS

Now let's look at some applications and how they fit in

Application 1: 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

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

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

$$v(w) = \max_{a \in \Gamma(x)} B(x, a, v) = \max \left\{ \frac{w}{1 - \beta}, c + \beta(R_\tau v)(w) \right\}$$

Proposition. \mathcal{R} is a contracting RDP

Proof: The quantile map R_τ obeys (see Ch. 7)

$$R_\tau(v + \lambda) = R_\tau v + \lambda \text{ for all } v \in \mathbb{R}^X \text{ and } \lambda \in \mathbb{R}$$

Hence

$$\begin{aligned} B(x, a, v + \lambda) &= a \frac{w}{1 - \beta} + (1 - a)[c + \beta(R_\tau(v + \lambda))(w)] \\ &= a \frac{w}{1 - \beta} + (1 - a)[c + \beta(R_\tau v)(w)] + (1 - a)\beta\lambda \\ &\leq B(x, a, v) + \beta\lambda \end{aligned}$$

Hence

$$\begin{aligned} B(x, a, v) &= B(x, a, v' + v - v') \\ &\leq B(x, a, v' + \|v - v'\|_\infty) \leq B(x, a, v') + \beta \|v - v'\|_\infty \end{aligned}$$

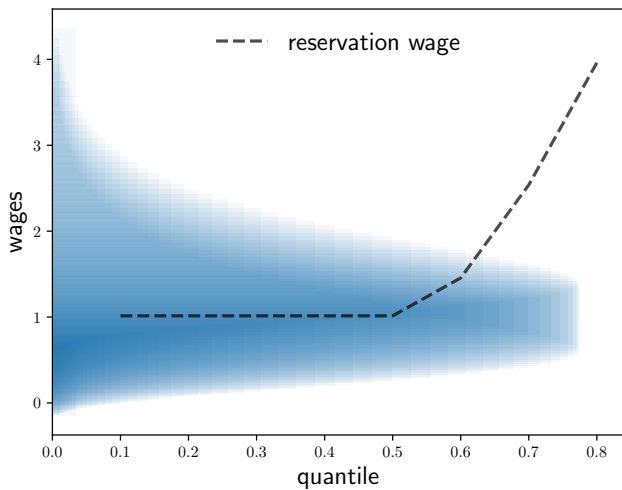
$$\therefore B(x, a, v) - B(x, a, v') \leq \beta \|v - v'\|_\infty$$

Reversing the roles of v and v' gives

$$|B(x, a, v) - B(x, a, v')| \leq \beta \|v - v'\|_\infty$$

Hence \mathcal{R} is contracting and, since V is closed, globally stable

Hence all optimality properties apply



Application 2: adversarial agents

Consider

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

Choice $d \in D$ is made by the adversary

$P(x, a, d, \cdot)$ is a distribution over X for each feasible (x, a, d)

We assume that

- Γ is a nonempty correspondence from X to A
- D is nonempty

Set

$$B(x, a, v) := \min_{d \in D} \left\{ r(x, a, d) + \beta \sum_{x'} v(x') P(x, a, d, x') \right\}$$

Fix $\varepsilon > 0$, set

$$V = [v_1, v_2] \quad \text{where} \quad v_1 := \frac{\min r - \varepsilon}{1 - \beta} \quad \text{and} \quad v_2 := \frac{\max r}{1 - \beta}$$

Then (Γ, V, B) is an RDP and

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

recovers the adversarial agent Bellman equation

Prop. \mathcal{R} is a concave RDP

Proof Fixing $(x, a) \in G$ and $v, w \in V$, we have

$$\begin{aligned} B(x, a, \lambda v + (1 - \lambda)w) &= \min_d \left\{ r + \beta \sum (\lambda v + (1 - \lambda)w)P \right\} \\ &= \min_d \left\{ \lambda [r + \beta \sum vP] + (1 - \lambda) [r + \beta \sum wP] \right\} \end{aligned}$$

Since $\min (f + g) \geq \min f + \min g$, we have

$$B(x, a, \lambda v + (1 - \lambda)w) \geq \lambda B(x, a, v) + (1 - \lambda)B(x, a, w)$$

For remaining minor details see Ch. 8

Application 3: inventory management

Consider an inventory problem with

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

where

- y is inventory, a is current order
- $r(y, a)$ is current profits
- $X := \{0, \dots, K\}$
- $\beta(z) = 1/(1 + r(z))$ — time-varying interest rates

This is an RDP $\mathcal{R} = (\Gamma, B, V)$ with

- $V := \mathbb{R}^X$ for $X := Y \times Z$
- $\Gamma(y, z) := \{0, \dots, K - y\}$ and

$$B((y, z), a, v) := r(y, a) + \beta(z) \sum_{z', y'} v(y', z') Q(z, z') R(y, a, y')$$

Proposition. If $L(z, z') := \beta(z)Q(z, z')$ obeys $\rho(L) < 1$, then \mathcal{R} is eventually contracting

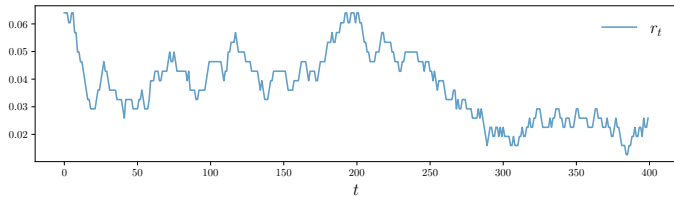
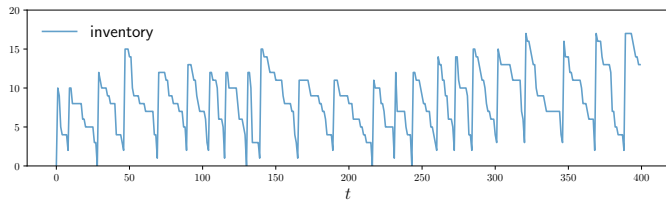
Proof: It suffices to show that, for all $(x, a) \in G$ and $v, w \in V$,

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

This holds because

$$\begin{aligned} & |B((y, z), a, v) - B((y, z), a, w)| \\ & \leq \beta(z) \sum_{y', z'} |v(y', z') - w(y', z')| Q(z, z') R(y, a, y') \\ & \leq \beta(z) \sum_{y', z'} |v(y', z') - w(y', z')| Q(z, z') \\ & = \sum_{y', z'} |v(y', z') - w(y', z')| L(z, z') \end{aligned}$$

Since $\rho(L) < 1$, we are done



Further applications

Other applications covered in the book:

- Epstein–Zin — details
- Robust control
- Smooth ambiguity, etc.
- Equilibria in production models
- Jump processes (continuous time)
- etc.