

Semicontractive Dynamic Programming

Dimitri P. Bertsekas

Department of Electrical Engineering and Computer Science
Massachusetts Institute of Technology

Lecture 4 of 5

September 2016

- Semicontractive Examples.
- Semicontractive Analysis for Stochastic Optimal Control.
- Extensions to Abstract DP Models.
- Applications to Stochastic Shortest Path and Other Problems.
- Algorithms.

Outline of this Lecture

- 1 Review of Abstract DP
- 2 Semicontractive Analysis
- 3 Stochastic Shortest Path Problem
- 4 Affine Monotonic Problem: Exponential Cost Function
- 5 Minimax Shortest Path Problem

Abstract DP Problem Formulation

- **State and control spaces:** X, U
- **Control constraint:** $u \in U(x)$ for all x
- **Stationary policies:** $\mu : X \mapsto U$, with $\mu(x) \in U(x)$ for all x

Monotone Mappings

- **Abstract monotone mapping** $H : X \times U \times E(X) \mapsto \mathfrak{R}$

$$J \leq J' \quad \implies \quad H(x, u, J) \leq H(x, u, J'), \quad \forall x, u$$

where $E(X)$ is the set of functions $J : X \mapsto [-\infty, \infty]$

- Mappings T_μ and T

$$(T_\mu J)(x) = H(x, \mu(x), J), \quad \forall x \in X, J \in E(X)$$

$$(TJ)(x) = \inf_{\mu} (T_\mu J)(x) = \inf_{u \in U(x)} H(x, u, J), \quad \forall x \in X, J \in E(X)$$

Stochastic Optimal Control Mapping: A Special Case

$$H(x, u, J) = E\{g(x, u, w) + \alpha J(f(x, u, w))\}$$

Abstract DP Problem

- Given an **initial function** $\bar{J} \in R(X)$ and policy μ , define

$$J_\mu(x) = \limsup_{N \rightarrow \infty} (T_\mu^N \bar{J})(x), \quad x \in X$$

- Find $J^*(x) = \inf_\mu J_\mu(x)$ and an optimal μ attaining the infimum

Results of Interest

- Bellman's equation**

$$J^* = TJ^*$$

and its set of solutions. Usually J^* is a solution.

- Conditions for optimality** of a stationary policy μ , usually $T_\mu J_\mu = TJ_\mu$.
- Algorithms, such as **value iteration (VI)** and **policy iteration (PI)**, and their convergence issues.

Semicontractive Models:

Some policies are "well-behaved" (have a regularity property), and others are not.

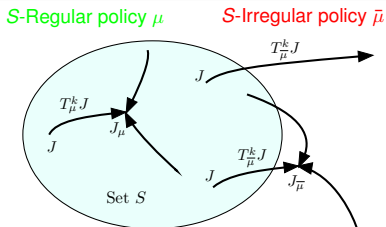
- Select a class of well-behaved/regular policies
- Define a **restricted optimization problem** over the regular policies only
- Show that the restricted problem has nice theoretical and algorithmic properties
- Relate the restricted problem to the original
- Under reasonable conditions, obtain strong theoretical and algorithmic results

Research Monograph

D. P. Bertsekas, Abstract Dynamic Programming, Athena Scientific, 2013; **updated chapters on-line**

S-Regularity

Key idea: We have a set of functions $S \subset E(X)$, which we view as the “domain of regularity”



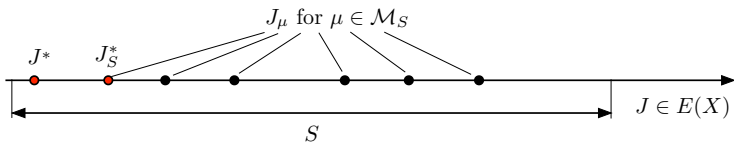
Definition of S-Regular Policy

Given a set of functions $S \subset E(X)$, we say that a stationary policy μ is **S-regular** if:

- $J_\mu \in S$ and $J_\mu = T_\mu J_\mu$
- $T_\mu^k J \rightarrow J_\mu$ for all $J \in S$

A policy that is not S-regular is called **S-irregular**.

S-Regular Restricted Problem



Given a set $S \subset E(X)$

- Consider the **restricted optimization problem**: Minimize J_μ over μ in the set \mathcal{M}_S of all S -regular policies
- Let J_S^* be the optimal cost function over S -regular policies only:

$$J_S^*(x) = \inf_{\mu \in \mathcal{M}_S} J_\mu(x), \quad x \in X$$

- Since the set of S -regular policies is a subset of the set of all policies,

$$J^* \leq J_S^*$$

A Principal Assumption that Guarantees “Good Behavior”

Assume that S consists of real-valued functions and:

- There exists at least one S -regular policy and $J_S^* = \inf_{\mu \in \mathcal{M}_S} J_\mu$ belongs to S .
- For every $J \in S$ and S -irregular policy μ , there exists $x \in X$ such that

$$\limsup_{k \rightarrow \infty} (T_\mu^k J)(x) = \infty$$

- S contains \bar{J} , and has the property that if J_1, J_2 are two functions in S , then S contains all functions J with $J_1 \leq J \leq J_2$
- The set $\{u \in U(x) \mid H(x, u, J) \leq \lambda\}$ is compact for every $J \in S$, $x \in X$, and $\lambda \in \mathfrak{R}$
- For each sequence $\{J_m\} \subset S$ with $J_m \uparrow J$ for some $J \in S$,

$$\lim_{m \rightarrow \infty} H(x, u, J_m) = H(x, u, J), \quad \forall x \in X, u \in U(x)$$

- For each function $J \in S$, there exists a function $J' \in S$ such that $J' \leq J$ and $J' \leq TJ'$

It is worth checking which parts of these assumptions are violated in the counterexamples of Lecture 1.

Proposition: Under the preceding assumption

- (Bellman Eq.) $J^* = TJ^*$. Moreover, J^* is the unique fixed point of T within S
- (VI Convergence) We have $T^k J \rightarrow J^*$ for all $J \in S$
- (Optimality Condition) μ is optimal if and only if $T_\mu J^* = TJ^*$, and there exists an optimal S -regular μ
- (PI Convergence) If in addition for each $\{J_m\} \subset E(X)$ with $J_m \downarrow J$ for some $J \in E(X)$,

$$H(x, u, J) = \lim_{m \rightarrow \infty} H(x, u, J_m), \quad \forall x \in X, u \in U(x)$$

then every sequence $\{\mu^k\}$ generated by the PI algorithm starting from an S -regular policy μ^0 satisfies $J_{\mu^k} \downarrow J^*$

- (Optimization-Based Solution of Bellman's Eq.) For any $J \in S$, if $J \leq TJ$ we have $J \leq J^*$, and if $J \geq TJ$ we have $J \geq J^*$

Note: Nearly as strong results as for contractive problems.

Common Characteristics

- They all involve a finite number of states ($X = \{1, \dots, n\}$), and a finite number of controls at each state (so **the number of policies is finite**).
- The set $R(X)$ of real-valued functions on X is identified with \mathfrak{R}^n .
- In all cases S is a subset of \mathfrak{R}^n .
- Usually there is a termination state.
- Because of this structure, **the complicated assumption given earlier simplifies**, and is nonrestrictive and intuitive.
- The results are almost as strong as for discounted problems.

- **Stochastic Shortest Path (SSP)** Problems: Transition probs. $p_{ij}(u)$,

$$\bar{J}(i) \equiv 0, \quad (T_\mu J)(i) = \sum_{j=1}^n p_{ij}(\mu(i)) (g(i, \mu(i), j) + J(j))$$

- **Affine Monotonic (AM)** Problems:

$$\bar{J} \geq 0, \quad T_\mu J = b_\mu + A_\mu J,$$

where $b_\mu \geq 0$, $A_\mu \geq 0$. A special case is **SSP with exponential cost**.

- **Minimax Shortest Path (MSP)** Problems: Disturbance has a nonprobabilistic set-membership description, $w \in W(i)$,

$$\bar{J}(i) \equiv 0, \quad (T_\mu J)(i) = \max_{w \in W(i)} \{g(i, \mu(i), w) + \alpha J(f(i, \mu(i), w))\}$$

$$J_\mu(i_0) = \limsup_{N \rightarrow \infty} \max_{w_0, w_1, \dots} \sum_{k=0}^N g(i_k, \mu(i_k), w_k)$$

We Specialize our Analysis to the Finite Spaces Context

S consists of real-valued functions

- For SSP and MSP, we use $S = \mathfrak{R}^n$.
- For AM, we use $S = \mathfrak{R}_+^n$, the nonnegative orthant.

Thanks to the finite spaces structure, and the choices of S , the complicated multipart assumption simplifies to the following:

- **There exists at least one S -regular policy.**
- **Infinite cost condition:** For all $J \in S$ and S -irregular μ , there exists i such that

$$\limsup_{k \rightarrow \infty} (T_\mu^k J)(i) = \infty$$

All other parts of the assumption are automatically satisfied.

A Common Approach for All Three Applications

Define the set S

For SSP and MSP, we use $S = \mathfrak{R}^n$. For AM, we use $S = \mathfrak{R}_+^n$, the nonnegative orthant.

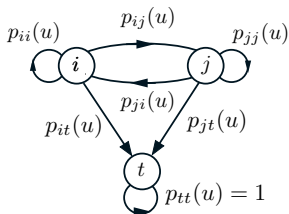
Characterize the S -Regular Policies

- For SSP, the S -regular μ are the **proper policies** (those that terminate with prob. 1).
- For AM, the S -regular μ are those for which **T_μ is a contraction**, i.e., all eigenvalues of A_μ are strictly within the unit circle.
- For MSP, the S -regular μ turn out to be those that **guarantee termination** regardless of the adversarial actions w_0, w_1, \dots , but **also some others**.

Assume that **there exists an S -regular policy** and that **each S -irregular policy has infinite cost**.

Apply the theorem: **J^* solves uniquely Bellman's Eq., VI, PI, and optimization approach work, etc.**

Stochastic Shortest Path Problem



A graph of n nodes plus the destination t

- At each node i we choose one of m probability distributions $p_{ij}(u)$, $u = 1, \dots, m$, over the successor nodes j .
- Transition cost $g(i, u, j)$.
- Minimize total expected cost up to termination.

$$\bar{J}(i) \equiv 0, \quad (T_{\mu}J)(i) = \sum_{i=1}^n p_{ij}(\mu(i)) (g(i, \mu(i), j) + J(j))$$

Proper Policies

- A policy μ is **proper** if it terminates from every initial state with probability 1.
- Equivalent definition: Starting at any node i , there exists a sequence of positive probability transitions under μ that starts at i and ends at t .
- Then $J_\mu(i)$ is the expected cost starting from i up to termination.

S-Regularity

- A policy is **S-regular**, where $S = \mathfrak{R}^n$, if and only if it is proper.
- We just verify the regularity definition ($T_\mu^k J \rightarrow J_\mu$ for all $J \in S$): We have that $T_\mu^k J$ does not depend on J for k large if and only if μ terminates.
- Assume there exists a proper policy.

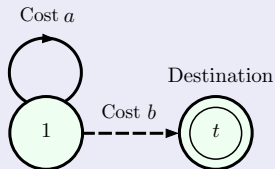
Assume each Improper Policy has Infinite Cost Starting at Some Initial State

Check that “cycling has positive cost”; true if every transition has positive cost.

Apply the theorem

J^* solves uniquely Bellman's Eq., VI and PI converge to J^* , LP approach works, etc.

Deterministic shortest path problem



One proper policy (from 1 go to t), and one improper policy (self-cycle)

Set of solutions of Bellman's equation: $J(1) = \min \{b, a + J(1)\}$

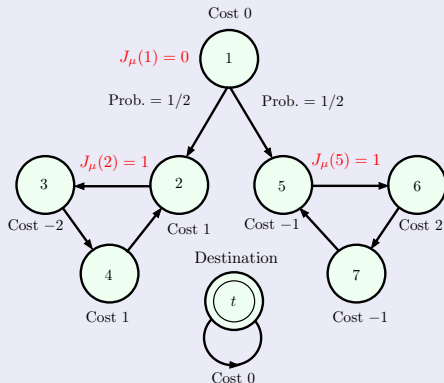
- Unique solution, $J^*(1) = b$ if $a > 0$ (assumptions satisfied)
- All $J(1) \leq b$ if $a = 0$ (assumptions violated)
- No real-valued solution if $a < 0$ (assumptions violated; consider changing S)

The assumption $a > 0$ corresponds to the classical conditions:

- There exists a path to the destination starting from every node.
- All cycles have positive length.

The SSP Problem where J^* does not Satisfy Bellman's Equation

A single policy μ . The only uncertainty is at the first stage starting at state 1



The Bellman Eq. is violated at 1: $J_\mu(1) \neq \frac{1}{2}J_\mu(2) + \frac{1}{2}J_\mu(5)$

Here the infinite cost condition is violated.

Affine Monotonic Problems

T_μ maps $J \in \mathfrak{R}_+^n$ into $T_\mu J \in \mathfrak{R}_+^n$ and is affine:

$$T_\mu J = b_\mu + A_\mu J,$$

where $b_\mu \geq 0$, $A_\mu \geq 0$. Also assume $\bar{J} \in \mathfrak{R}_+^n$ (but may have $\bar{J} \neq 0$)

Some special cases

- An **SSP problem with nonnegative cost** per transition. Corresponds to $\bar{J} = 0$ and

$$b_\mu(i) = g(i, \mu(i), j), \quad A_\mu(i, j) = p_{ij}(\mu(i))$$

- An SSP problem with **exponential cost** for the length of a path, so

$$J_\mu(i) = E\{\exp(\text{Length of path starting at } i \text{ up to reaching destination } t)\}$$

Corresponds to the affine monotonic problem defined by

$$\bar{J}(i) \equiv 1, \quad (T_\mu J)(i) = p_{it}(\mu(i)) e^{g(i, \mu(i), t)} + \sum_{j=1}^n p_{ij}(\mu(i)) e^{g(i, \mu(i), j)} J(j)$$

- **Multiplicative cost function** (contains the exponential cost SSP as a special case)

Cost Function of a Policy μ

By repeatedly applying the equation $T_\mu J = b_\mu + A_\mu J$, we have

$$T_\mu^N J = A_\mu^N J + \sum_{k=0}^{N-1} A_\mu^k b_\mu, \quad \forall J \in E^+(X), N = 1, 2, \dots,$$

$$J_\mu = \limsup_{N \rightarrow \infty} T_\mu^N \bar{J} = \limsup_{N \rightarrow \infty} A_\mu^N \bar{J} + \sum_{k=0}^{\infty} A_\mu^k b_\mu$$

Contractive policies: Those for which $\limsup_{N \rightarrow \infty} A_\mu^N J = 0$ for all $J \in \mathfrak{R}^n$ (equivalently A_μ has eigenvalues strictly within the unit circle).

Key fact is that μ is $R^+(X)$ -regular if and only if T_μ is contractive. Justification:

$$J_\mu = \limsup_{N \rightarrow \infty} T_\mu^N J = \limsup_{N \rightarrow \infty} \sum_{k=0}^{N-1} A_\mu^k b_\mu, \quad \forall \mu: \text{contractive}, J \in \mathfrak{R}_+^n$$

Hence, if μ is contractive it is also $R^+(X)$ -regular. The reverse can also be shown to be true.

Assume that:

- There exists at least one contractive policy
- Each noncontractive policy has infinite cost for some initial state.

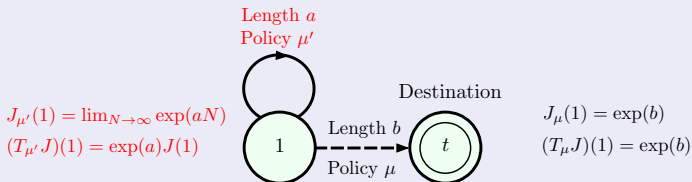
Then the standard results hold:

- Bellman's Eq. has J^* as its unique solution
- VI and PI converge to J^*
- Standard optimality conditions hold
- Solution by linear programming is possible

Some notes for the exponential cost SSP

- **Every proper policy is contractive but the reverse is not true** (consider a deterministic problem and a policy with a negative length cycle)
- **In exponential cost SSP policies that include cycles with "negative cost" do not cause difficulties** (but "zero cost cycles" may cause a problem)

Back to the Deterministic Shortest Path Problem



Bellman's equation: $J(1) = \min \{ \exp(b), \exp(a)J(1) \}$

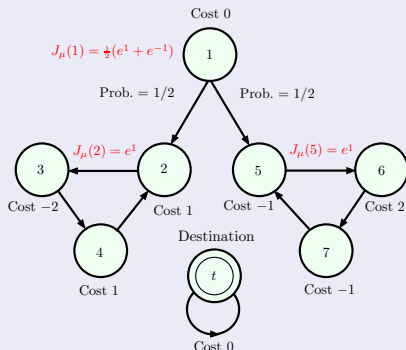
- If $a > 0$ (assumptions satisfied), $J^*(1) = \exp(b)$ solves uniquely Bellman's Eq., μ is optimal
- If $a < 0$ (assumptions satisfied, **both policies are contractive, even though μ' is improper!**), $J^*(1) = 0$ solves uniquely Bellman's Eq., μ' is optimal
- If $a = 0$ (assumptions violated), all $J(1)$ in the interval $0 \leq J(1) \leq \exp(b)$ solve Bellman's Eq., $J^*(1) = \min \{ \exp(b), 1 \}$

The assumption $a > 0$ corresponds to the classical conditions:

- There exists a path to the destination starting from every node.
- All cycles have positive length.

An Exponential Cost SSP Problem where J^* does not Satisfy Bellman's Equation

This is the exponential cost version of the earlier SSP counterexample, which involved zero length cycles.



The Bellman Eq. is violated at 1: $J_\mu(1) \neq \frac{1}{2}J_\mu(2) + \frac{1}{2}J_\mu(5)$

Here the policy is noncontractive and hence \mathfrak{R}_+^n -irregular, while the infinite cost condition is violated.

Problem Formulation

- A graph with set of nodes $X = \{1, \dots, n\}$ plus a destination t , and a set of directed arcs (i, j) , where $i, j \in X \cup \{t\}$.
- At each node i we may choose a control u from a finite set $U(i)$.
- The destination t is absorbing and cost-free.
- At node i , a successor node j is selected by an antagonistic opponent from a given set $Y(i, u) \subset X \cup \{t\}$ and a cost $g(i, u, j)$ is incurred.
- Mappings:

$$H(i, u, J) = \max_{j \in Y(i, u)} [g(i, u, j) + \tilde{J}(j)], \quad \forall x, u, J \in \mathbb{R}^n,$$

where $\tilde{J}(j) = J(j)$ if $j \in X$ and $\tilde{J}(j) = 0$ if $j = t$. We have

$$(T_\mu J)(i) = H(i, \mu(i), J), \quad (TJ)(i) = \min_{u \in U(i)} H(i, u, J)$$

- Let \bar{J} be the zero function, so

$$J_\mu(i_0) = \limsup_{N \rightarrow \infty} \max_{w_0, w_1, \dots} \sum_{k=0}^N g(i_k, \mu(i_k), w_k)$$

Cost Function and Other Properties of a Policy μ

- A **possible path under μ** starting at node $i_0 \in X$ is an arc sequence $p = \{(i_0, i_1), (i_1, i_2), \dots\}$, such that $i_{k+1} \in Y(i_k, \mu(i_k))$ for all $k \geq 0$. The set of all possible paths under μ starting at i_0 is denoted by $P(i_0, \mu)$.
- The **length of a path** $p \in P(i_0, \mu)$ is $\limsup_{N \rightarrow \infty} \sum_{k=0}^N g(i_k, \mu(i_k), i_{k+1})$.
- Similar definitions for the length of a **portion of a path** p , consisting of a finite number of consecutive arcs.
- For any μ and i , $(T_\mu^k \bar{J})(i)$ is the **length of the longest path under μ that starts at i and consists of k arcs**, and can be computed with a k -stage DP algorithm.
- Of special interest are **cycles**, i.e., paths of the form $\{(i_j, i_{j+1}), \dots, (i_{j+m}, i_j)\}$, and paths that **terminate**, i.e., have the form $p = \{(i_0, i_1), \dots, (i_m, t), (t, t), \dots\}$.

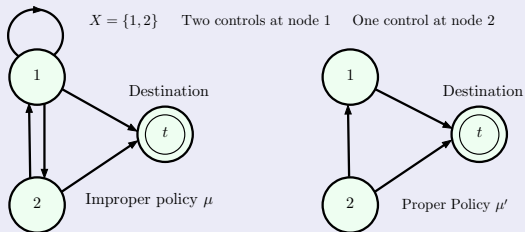
Proper Policies

A policy μ is **proper** if for all i , all the paths in $P(i, \mu)$ contain no cycle and terminate.

S-Regular Policies ($S = \mathfrak{R}^n$)

It is easy to see that **all proper policies are \mathfrak{R}^n -regular**. The reverse is not true.

The Characteristic Graph of a Policy μ : $\mathcal{A}_\mu = \cup_{i \in X} \{(i, j) \mid j \in Y(i, \mu(i))\}$



- We say that \mathcal{A}_μ is **destination-connected** if for each $i \in X$ there exists a terminating path in $P(i, \mu)$.

Characterization of \mathfrak{R}^n -Regular Policies

- μ is \mathfrak{R}^n -regular if and only if \mathcal{A}_μ is destination-connected and all its cycles have negative length. (Note that **a proper policy is \mathfrak{R}^n -regular.**)
- μ is \mathfrak{R}^n -irregular if and only if it is improper, and either is destination-disconnected or \mathcal{A}_μ has a cycle with length ≥ 0 . (Note that **there exist improper policies that are \mathfrak{R}^n -regular.**)

Assume that:

- **There exists at least one proper policy** (implies that there exists an \mathfrak{R}^n -regular policy).
- **For every improper policy μ , all cycles in the characteristic graph \mathcal{A}_μ have positive length** (implies that every \mathfrak{R}^n -irregular policy has infinite cost for some initial state).

Then the standard results hold:

- Bellman's Eq. has J^* as its unique solution.
- VI, PI, converge to J^* .
- Standard optimality conditions hold, etc.

Some notes

- **The positive cycle condition can be relaxed to nonnegativity**, using a perturbation approach (add a $\delta > 0$ to each $g(i, u, j)$ and take $\delta \downarrow 0$; see the next lecture).
- **There is a finitely terminating Dijkstra-like algorithm** for MSP problems with nonnegative arc lengths (this is a consequence of the shortest path character of the problem, not its semicontractive character).