

The Value Function Polytope in Reinforcement Learning

Robert Dadashi, Adrien Ali Taiga, Nicolas Le Roux, Dale Schuurmans, Marc G. Bellemare

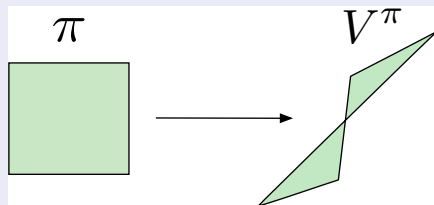
Characterizing the space of value functions in a finite state-action Markov Decision Process context

Characterizing the space of value functions in a finite
state-action Markov Decision Process context

A geometric perspective

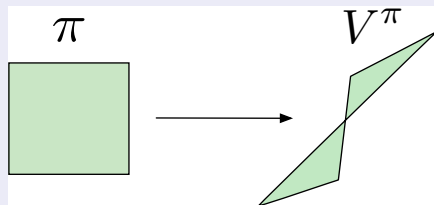
Main Idea

The space of value functions forms a polytope

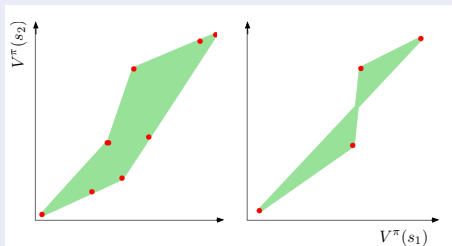


Main Idea

The space of value functions forms a polytope



The boundary of this polytope can be described by value functions corresponding to “semi-deterministic” policies



Preliminaries

We are working with a Markov Decision Process, $M := (\mathcal{S}, \mathcal{A}, r, P, \gamma)$

\mathcal{S} : finite state space

\mathcal{A} : finite action space

r : reward function

P : transition function

γ : discount factor in $[0, 1)$

Preliminaries

We are working with a Markov Decision Process, $M := (\mathcal{S}, \mathcal{A}, r, P, \gamma)$

\mathcal{S} : finite state space

\mathcal{A} : finite action space

r : reward function

P : transition function

γ : discount factor in $[0, 1)$

π denotes a policy; $\pi : \mathcal{S} \rightarrow \Delta\mathcal{A}$. Together with the transition function P , we have state-to-state transition probabilities with respect to policy π :

Preliminaries

We are working with a Markov Decision Process, $M := (\mathcal{S}, \mathcal{A}, r, P, \gamma)$

\mathcal{S} : finite state space

\mathcal{A} : finite action space

r : reward function

P : transition function

γ : discount factor in $[0, 1)$

π denotes a policy; $\pi : \mathcal{S} \rightarrow \Delta\mathcal{A}$. Together with the transition function P , we have state-to-state transition probabilities with respect to policy π :

$$P^\pi(s'|s) := \sum_{a \in \mathcal{A}} \pi(a|s)P(s'|s, a)$$

The value function at state s is defined as follows:

$$V^\pi(s) := \mathbb{E}_{P^\pi} \left[\sum_{i=0}^{\infty} \gamma^i r(s_i, a_i) \mid s_0 = s \right]$$

where $r(s, a)$ is the reward function evaluated at state s and action a

The value function at state s is defined as follows:

$$V^\pi(s) := \mathbb{E}_{P^\pi} \left[\sum_{i=0}^{\infty} \gamma^i r(s_i, a_i) \mid s_0 = s \right]$$

where $r(s, a)$ is the reward function evaluated at state s and action a

Definition

Let $f_v : \mathcal{P}(\mathcal{A})^S \rightarrow \mathbb{R}^S$ be the **value functional** mapping the space of policies to their corresponding value functions.

Definition

Policy Determinism: A policy π is

- **s-deterministic** for $s \in \mathcal{S}$ if $\pi(a|s) \in \{0, 1\}$.
- **semi-deterministic** if it is s-deterministic for at least one $s \in \mathcal{S}$.
- **deterministic** if it is s-deterministic for all states $s \in \mathcal{S}$.

Definition

Policy Determinism: A policy π is

- **s-deterministic** for $s \in \mathcal{S}$ if $\pi(a|s) \in \{0, 1\}$.
- **semi-deterministic** if it is s -deterministic for at least one $s \in \mathcal{S}$.
- **deterministic** if it is s -deterministic for all states $s \in \mathcal{S}$.

Definition

Let s_1, \dots, s_k be states and π a policy. Then $\mathbf{Y}_{s_1, \dots, s_k}^\pi \subset \mathcal{P}(\mathcal{A})^{\mathcal{S}}$ is the set of policies that agree with π on the states s_1, \dots, s_k . Similarly, let $\mathbf{Y}_{\mathcal{S}-\{s\}}^\pi$ be the set of policies which agree with π on all states except s .

Visualizations

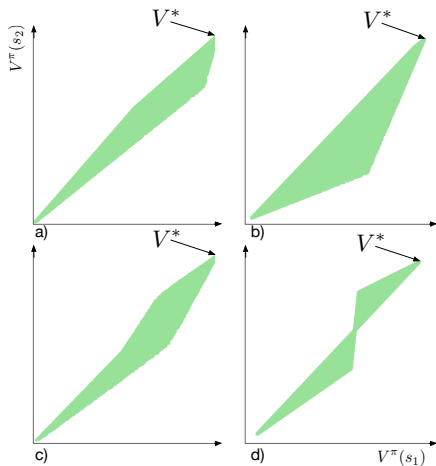


Figure: Value function space corresponding to four two-state MDPs; each evaluated from 50,000 policies sampled uniformly from $\mathcal{P}(\mathcal{A})^S$.

The Line Theorem

Theorem 1: Let s be a state and π a policy. Then there are two s -deterministic policies in $Y_{S-\{s\}}^\pi$, denoted π_l and π_u , such that for all $\pi' \in Y_{S-\{s\}}^\pi$

$$f_v(\pi_l) \preceq f_v(\pi') \preceq f_v(\pi_u)$$

Furthermore, the following are equivalent:

- $f_v(Y_{S-\{s\}}^\pi)$
- $\{f_v(\alpha\pi_l + (1-\alpha)\pi_u) : \alpha \in [0, 1]\}$
- $\{\alpha f_v(\pi_l) + (1-\alpha)f_v(\pi_u) : \alpha \in [0, 1]\}$

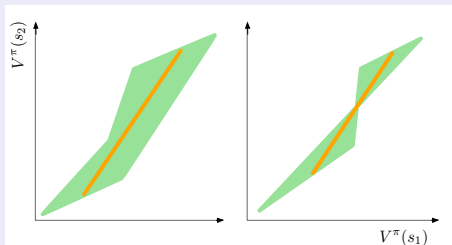
The Line Theorem

Theorem 1: Let s be a state and π a policy. Then there are two s -deterministic policies in $Y_{S-\{s\}}^\pi$, denoted π_l and π_u , such that for all $\pi' \in Y_{S-\{s\}}^\pi$

$$f_v(\pi_l) \preceq f_v(\pi') \preceq f_v(\pi_u)$$

Furthermore, the following are equivalent:

- $f_v(Y_{S-\{s\}}^\pi)$
- $\{f_v(\alpha\pi_l + (1-\alpha)\pi_u) : \alpha \in [0, 1]\}$
- $\{\alpha f_v(\pi_l) + (1-\alpha)f_v(\pi_u) : \alpha \in [0, 1]\}$



Definitions

Definition

A **convex combination** is a finite linear combination of vectors whose coefficients are non-negative and sum to 1.

Definitions

Definition

A **convex combination** is a finite linear combination of vectors whose coefficients are non-negative and sum to 1.

Definition

The **convex hull** of a set E is the intersection of all convex sets containing E . A set C is **convex** if for any $x, y \in C$
 $\{\alpha x + (1 - \alpha)y : \alpha \in [0, 1]\} \subset C$.

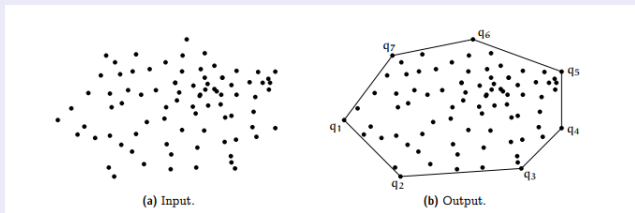
Definitions

Definition

A **convex combination** is a finite linear combination of vectors whose coefficients are non-negative and sum to 1.

Definition

The **convex hull** of a set E is the intersection of all convex sets containing E . A set C is **convex** if for any $x, y \in C$
 $\{\alpha x + (1 - \alpha)y : \alpha \in [0, 1]\} \subset C$.



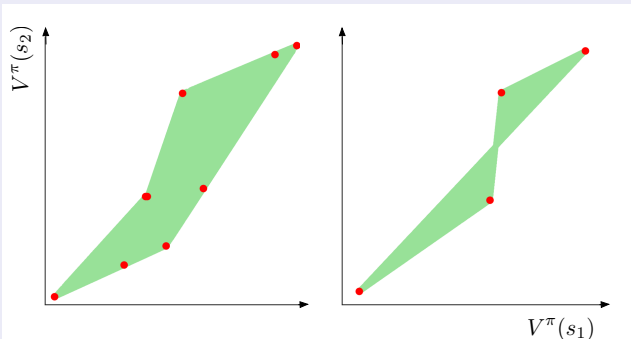
Credit to Harshit Sikchi.

Interesting Consequences of the Line Theorem

For any set of states $s_1, \dots, s_k \in \mathcal{S}$ and a policy π , V^π can be expressed as a convex combination of value functions of $\{s_1, \dots, s_k\}$ -deterministic policies. In particular, $\mathcal{V} := f_v(\mathcal{P}(\mathcal{A})^{\mathcal{S}})$ is included in the convex hull of the value functions of deterministic policies.

Interesting Consequences of the Line Theorem

For any set of states $s_1, \dots, s_k \in \mathcal{S}$ and a policy π , V^π can be expressed as a convex combination of value functions of $\{s_1, \dots, s_k\}$ -deterministic policies. In particular, $\mathcal{V} := f_v(\mathcal{P}(\mathcal{A})^{\mathcal{S}})$ is included in the convex hull of the value functions of deterministic policies.



Interesting Consequences of the Line Theorem

Let V^π and $V^{\pi'}$ be two value functions. Then there exists a sequence of policies π_1, \dots, π_k ($k \leq \mathcal{S}$ such that $V^\pi = V^{\pi_1}$, $V^{\pi'} = V^{\pi_k}$, and

$$\{f_v(\alpha\pi_i + (1 - \alpha)\pi_{i+1}) : \alpha \in [0, 1]\}$$

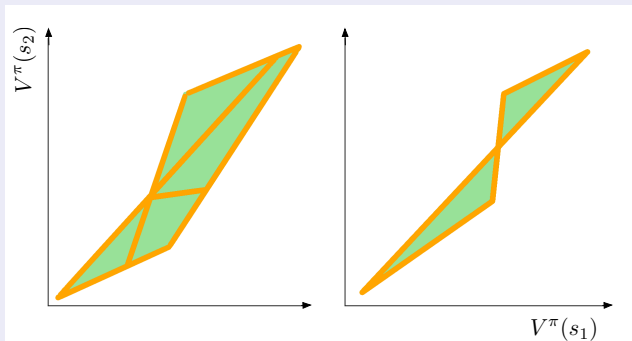
forms a line segment for all $1 \leq i < k$.

Boundary of Semi-Deterministic Policies

The boundary of the space of value functions is a subset of value functions corresponding to semi-deterministic policies.

Boundary of Semi-Deterministic Policies

The boundary of the space of value functions is a subset of value functions corresponding to semi-deterministic policies.



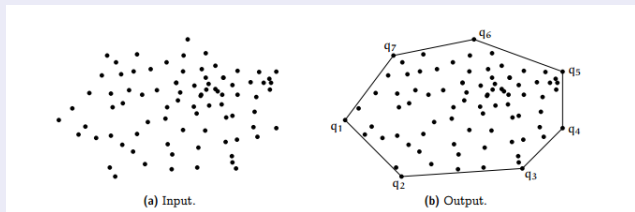
Definitions

Definition

P is a **convex polytope** iff there exist points $x_1, \dots, x_k \in \mathbb{R}^n$ such that P is the convex hull of $\{x_1, \dots, x_k\}$.

Definition

A **polytope** is a finite union of convex polytopes.



Credit to Harshit Sikchi.

Main Result

Let π be a policy and let s_1, \dots, s_k be states in \mathcal{S} . Then $f_v(Y_{s_1, \dots, s_k}^\pi)$ is a polytope and in particular, $\mathcal{V} = f_v(Y_\phi^\pi)$ is a polytope.

Main Result

Let π be a policy and let s_1, \dots, s_k be states in \mathcal{S} . Then $f_v(Y_{s_1, \dots, s_k}^\pi)$ is a polytope and in particular, $\mathcal{V} = f_v(Y_\phi^\pi)$ is a polytope.

- Surprising since f_v is in general non-linear and mixtures of policies can describe curves
- There is a sub-polytope structure in the space of value functions

