Understanding (Exact) Dynamic Programming through Bellman Operators

Ashwin Rao

ICME, Stanford University

January 14, 2019

Overview

- Vector Space of Value Functions
- 2 Bellman Operators
- 3 Contraction and Monotonicity
- Policy Evaluation
- 6 Policy Iteration
- **6** Value Iteration
- Policy Optimality

Vector Space of Value Functions

- Assume State pace S consists of n states: $\{s_1, s_2, \ldots, s_n\}$
- Assume Action space A consists of m actions $\{a_1, a_2, \ldots, a_m\}$
- This exposition extends easily to continuous state/action spaces too
- We denote a stochastic policy as $\pi(a|s)$ (probability of "a given s")
- Abusing notation, deterministic policy denoted as $\pi(s) = a$
- ullet Consider a *n*-dim vector space, each dim corresponding to a state in ${\cal S}$
- \bullet A vector in this space is a specific Value Function (VF) $\textbf{v} \colon\thinspace \mathcal{S} \to \mathbb{R}$
- With coordinates $[\mathbf{v}(s_1), \mathbf{v}(s_2), \dots, \mathbf{v}(s_n)]$
- ullet Value Function (VF) for a policy π is denoted as $oldsymbol{v}_\pi:\mathcal{S} o\mathbb{R}$
- ullet Optimal VF denoted as $oldsymbol{v}_*: \mathcal{S}
 ightarrow \mathbb{R}$ such that for any $s \in \mathcal{S}$,

$$\mathbf{v}_*(s) = \max_{\pi} \mathbf{v}_{\pi}(s)$$



Some more notation

- ullet Denote \mathcal{R}^a_s as the Expected Reward upon action a in state s
- ullet Denote $\mathcal{P}^a_{s,s'}$ as the probability of transition s o s' upon action a
- Define

$$\mathsf{R}_{\pi}(s) = \sum_{\mathsf{a} \in \mathcal{A}} \pi(\mathsf{a}|s) \cdot \mathcal{R}_{\mathsf{s}}^{\mathsf{a}}$$

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

- Denote \mathbf{R}_{π} as the vector $[\mathbf{R}_{\pi}(s_1),\mathbf{R}_{\pi}(s_2),\ldots,\mathbf{R}_{\pi}(s_n)]$
- Denote \mathbf{P}_{π} as the matrix $[\mathbf{P}_{\pi}(s_i,s_{i'})], 1 \leq i,i' \leq n$
- ullet Denote γ as the MDP discount factor

Bellman Operators ${\bf B}_{\pi}$ and ${\bf B}_*$

- We define operators that transform a VF vector to another VF vector
- Bellman Policy Operator \mathbf{B}_{π} (for policy π) operating on VF vector \mathbf{v} :

$$\mathbf{B}_{\pi}\mathbf{v} = \mathbf{R}_{\pi} + \gamma \mathbf{P}_{\pi} \cdot \mathbf{v}$$

- ${\bf B}_{\pi}$ is a linear operator with fixed point ${\bf v}_{\pi}$, meaning ${\bf B}_{\pi}{\bf v}_{\pi}={\bf v}_{\pi}$
- Bellman Optimality Operator \mathbf{B}_* operating on VF vector \mathbf{v} :

$$(\mathbf{B}_*\mathbf{v})(s) = \max_{a} \{\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{s,s'}^a \cdot \mathbf{v}(s')\}$$

- \mathbf{B}_* is a non-linear operator with fixed point \mathbf{v}_* , meaning $\mathbf{B}_*\mathbf{v}_* = \mathbf{v}_*$
- Define a function G mapping a VF v to a deterministic "greedy" policy $G(\mathbf{v})$ as follows:

$$G(\mathbf{v})(s) = \arg\max_{a} \{\mathcal{R}_{s}^{a} + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{s,s'}^{a} \cdot \mathbf{v}(s')\}$$

• $\mathbf{B}_{G(\mathbf{v})}\mathbf{v} = \mathbf{B}_*\mathbf{v}$ for any VF \mathbf{v} (Policy $G(\mathbf{v})$ achieves the max in \mathbf{B}_*)

Contraction and Monotonicity of Operators

- Both ${\bf B}_{\pi}$ and ${\bf B}_{*}$ are γ -contraction operators in L^{∞} norm, meaning:
- For any two VFs $\mathbf{v_1}$ and $\mathbf{v_2}$,

$$\|\mathbf{B}_{\pi}\mathbf{v}_{1}-\mathbf{B}_{\pi}\mathbf{v}_{2}\|_{\infty}\leq\gamma\|\mathbf{v}_{1}-\mathbf{v}_{2}\|_{\infty}$$

$$\|\mathbf{B}_*\mathbf{v_1} - \mathbf{B}_*\mathbf{v_2}\|_{\infty} \leq \gamma \|\mathbf{v_1} - \mathbf{v_2}\|_{\infty}$$

- So we can invoke Contraction Mapping Theorem to claim fixed point
- We use the notation $\mathbf{v_1} \leq \mathbf{v_2}$ for any two VFs $\mathbf{v_1}, \mathbf{v_2}$ to mean:

$$\mathsf{v_1}(s) \leq \mathsf{v_2}(s)$$
 for all $s \in \mathcal{S}$

- Also, both \mathbf{B}_{π} and \mathbf{B}_{*} are monotonic, meaning:
- For any two VFs $\mathbf{v_1}$ and $\mathbf{v_2}$,

$$v_1 \leq v_2 \Rightarrow B_{\pi}v_1 \leq B_{\pi}v_2$$

$$v_1 \leq v_2 \Rightarrow B_*v_1 \leq B_*v_2$$

Policy Evaluation

- ullet ${f B}_{\pi}$ satisfies the conditions of Contraction Mapping Theorem
- ullet ${f B}_{\pi}$ has a unique fixed point ${f v}_{\pi}$, meaning ${f B}_{\pi}{f v}_{\pi}={f v}_{\pi}$
- This is a succinct representation of Bellman Expectation Equation
- ullet Starting with any VF $oldsymbol{v}$ and repeatedly applying $oldsymbol{B}_{\pi}$, we will reach $oldsymbol{v}_{\pi}$

$$\lim_{N o \infty} \mathbf{B}_{\pi}^N \mathbf{v} = \mathbf{v}_{\pi}$$
 for any VF \mathbf{v}

This is a succinct representation of the Policy Evaluation Algorithm

Policy Improvement

- Let π_k and \mathbf{v}_{π_k} denote the Policy and the VF for the Policy in iteration k of Policy Iteration
- Policy Improvement Step is: $\pi_{k+1} = G(\mathbf{v}_{\pi_k})$, i.e. deterministic greedy
- Earlier we argued that $B_*v = B_{G(v)}v$ for any VF v. Therefore,

$$\mathbf{B}_* \mathbf{v}_{\pi_k} = \mathbf{B}_{G(\mathbf{v}_{\pi_k})} \mathbf{v}_{\pi_k} = \mathbf{B}_{\pi_{k+1}} \mathbf{v}_{\pi_k} \tag{1}$$

• We also know from operator definitions that $\mathbf{B}_*\mathbf{v} \geq \mathbf{B}_{\pi}\mathbf{v}$ for all π, \mathbf{v}

$$\mathbf{B}_* \mathbf{v}_{\pi_{\mathbf{k}}} \ge \mathbf{B}_{\pi_k} \mathbf{v}_{\pi_{\mathbf{k}}} = \mathbf{v}_{\pi_{\mathbf{k}}} \tag{2}$$

• Combining (1) and (2), we get:

$$\mathbf{B}_{\pi_{k+1}}\mathbf{v}_{\pi_{\mathbf{k}}} \geq \mathbf{v}_{\pi_{\mathbf{k}}}$$

• Monotonicity of $\mathbf{B}_{\pi_{k+1}}$ implies

$$\mathbf{B}_{\pi_{k+1}}^{\mathcal{N}}\mathbf{v}_{\pi_{\mathbf{k}}} \geq \ldots \mathbf{B}_{\pi_{k+1}}^{2}\mathbf{v}_{\pi_{\mathbf{k}}} \geq \mathbf{B}_{\pi_{k+1}}\mathbf{v}_{\pi_{\mathbf{k}}} \geq \mathbf{v}_{\pi_{\mathbf{k}}}$$
 $\mathbf{v}_{\pi_{\mathbf{k}+1}} = \lim_{N o \infty} \mathbf{B}_{\pi_{k+1}}^{N}\mathbf{v}_{\pi_{\mathbf{k}}} \geq \mathbf{v}_{\pi_{\mathbf{k}}}$

Policy Iteration

- ullet We have shown that in iteration k+1 of Policy Iteration, ${f v}_{\pi_{{f k}+1}} \geq {f v}_{\pi_{{f k}}}$
- If $\mathbf{v}_{\pi_{\mathbf{k}+1}} = \mathbf{v}_{\pi_{\mathbf{k}}}$, the above inequalities would hold as equalities
- ullet So this would mean $oldsymbol{\mathsf{B}}_*oldsymbol{\mathsf{v}}_{\pi_{oldsymbol{\mathsf{k}}}}=oldsymbol{\mathsf{v}}_{\pi_{oldsymbol{\mathsf{k}}}}$
- But B_{*} has a unique fixed point v_{*}
- ullet So this would mean $oldsymbol{v}_{\pi_{oldsymbol{k}}} = oldsymbol{v}_*$
- \bullet Thus, at each iteration, Policy Iteration either strictly improves the VF or achieves the optimal VF \textbf{v}_*

Value Iteration

- B_{*} satisfies the conditions of Contraction Mapping Theorem
- ullet $oldsymbol{\mathsf{B}}_*$ has a unique fixed point $oldsymbol{\mathsf{v}}_*$, meaning $oldsymbol{\mathsf{B}}_*oldsymbol{\mathsf{v}}_*=oldsymbol{\mathsf{v}}_*$
- This is a succinct representation of Bellman Optimality Equation
- \bullet Starting with any VF \boldsymbol{v} and repeatedly applying $\boldsymbol{B}_*,$ we will reach \boldsymbol{v}_*

$$\lim_{N o\infty} \mathbf{B}_*^N \mathbf{v} = \mathbf{v}_*$$
 for any VF \mathbf{v}

This is a succinct representation of the Value Iteration Algorithm

Greedy Policy from Optimal VF is an Optimal Policy

ullet Earlier we argued that $oldsymbol{B}_{G(oldsymbol{v})}oldsymbol{v} = oldsymbol{B}_*oldsymbol{v}$ for any VF $oldsymbol{v}$. Therefore,

$$\mathsf{B}_{G(\mathsf{v}_*)}\mathsf{v}_*=\mathsf{B}_*\mathsf{v}_*$$

• But \mathbf{v}_* is the fixed point of \mathbf{B}_* , meaning $\mathbf{B}_*\mathbf{v}_*=\mathbf{v}_*$. Therefore,

$$\mathsf{B}_{G(\mathsf{v}_*)}\mathsf{v}_*=\mathsf{v}_*$$

ullet But we know that ${f B}_{G({f v}_*)}$ has a unique fixed point ${f v}_{G({f v}_*)}$. Therefore,

$$\mathbf{v}_* = \mathbf{v}_{G(\mathbf{v}_*)}$$

- This says that simply following the deterministic greedy policy $G(\mathbf{v}_*)$ (created from the Optimal VF \mathbf{v}_*) in fact achieves the Optimal VF \mathbf{v}_*
- ullet In other words, $G(v_*)$ is an Optimal (Deterministic) Policy