
Rank-One Modified Value Iteration

Arman S. Kolarijani¹

Tolga Ok¹

Peyman Mohajerin Esfahani^{1, 2}

Mohamad Amin Sharif Kolarijani¹

¹Delft University of Technology, **DCSC**,
The Netherlands

²University of Toronto, **MIE**,
Canada



Generalized Policy Iteration

$$v_{k+1} = v_k + \mathbf{G}_k (\mathbf{T}(v_k) - v_k)$$

Policy Iteration

$$\mathbf{G}_k = (I - \gamma P_k)^{-1}$$

Value Iteration

$$\mathbf{G}_k = I$$

\mathbf{T} : Bellman optimality operator

π_v : Greedy policy w.r.t v

v_k : Value function at the k -th step

$P_k := P^{\pi_{v_k}}$: Transition matrix of the greedy policy π_{v_k}

Generalized Policy Iteration

$$v_{k+1} = v_k + \mathbf{G}_k (\mathbf{T}(v_k) - v_k)$$

Policy Iteration

$$\mathbf{G}_k = (I - \gamma P_k)^{-1}$$

Modified Policy Iteration

$$\mathbf{G}_k = \sum_{i=0}^{n-1} (\gamma P_k)^i$$

Value Iteration (VI)

$$\mathbf{G}_k = I$$

λ -Policy Iteration

$$\mathbf{G}_k = (I - \lambda \gamma P_k)^{-1}$$

Acceleration Methods

Planning	Learning
$v \in \mathbb{R}^n$	$q \in \mathbb{R}^{nm}$
$\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^n$	$[\widehat{\mathbf{T}}(q, \hat{s}^+)](s, a)$
$P^\pi \in \mathbb{R}^{n \times n}$	$\hat{s}^+ \sim P(\cdot s, a), \forall (s, a) \in \mathcal{S} \times \mathcal{A}$
→ Anderson VI	
→ Nesterov VI	Speedy Q-learning
→	Zap Q-learning
→ Operator Splitting VI	
→ Rank-One VI (R1VI)	Rank-One QL (R1QL)

Rank-One Value Iteration (R1VI)

Policy Iteration

$$v_{k+1} = v_k + \mathbf{G}_k (\mathbf{T}(v_k) - v_k)$$

Approximate P_k

$$\mathbf{G}_k = (I - \gamma P_k)^{-1}$$

$$\mathbf{G}_k = (I - \gamma \tilde{P}_k)^{-1}$$

Lemma 1: If the MDP is ergodic, then

$$\tilde{P}_k := \mathbf{1} \mathbf{d}_k^\top = \arg \min_{P \in \mathbb{R}^{n \times n}} \rho(P - P_k)$$

$$\text{s.t. } P \geq 0, \quad P\mathbf{1} = \mathbf{1}, \quad \text{rank}(P) = 1$$

Spectral radius

$$\mathbf{G}_k = (I - \gamma \tilde{P}_k)^{-1} \xrightarrow{\text{Rank-one approximation}} (I - \gamma \mathbf{1} \mathbf{d}_k^\top)^{-1} \xrightarrow{\text{Woodbury's formula}} I + \frac{\gamma}{1 - \gamma} \mathbf{1} \mathbf{d}_k^\top$$

$$v_{k+1} = \mathbf{T}^{\mathbf{R1}}(v_k) = \mathbf{T}(v_k) + \frac{\gamma}{1 - \gamma} \langle \mathbf{d}_k^\top, \mathbf{T}(v_k) - v_k \rangle \mathbf{1}$$

Rank-One Value Iteration (R1VI)

Algorithm R1VI

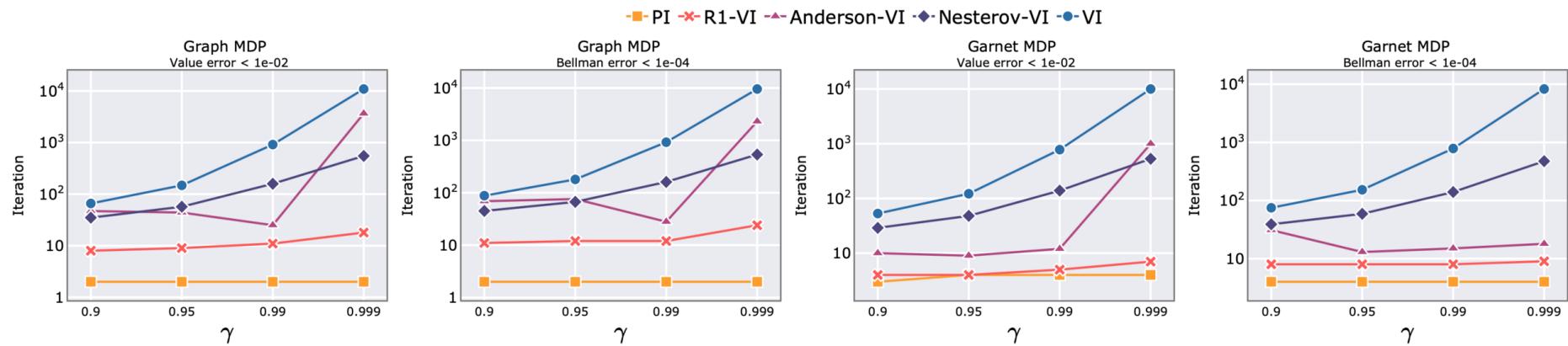
- **initialize** : $v_0 \in \mathbb{R}^n$, $\mathbf{d}_{-1} \in \Delta(\mathbb{R}^n)$
 - for $k = 0, 1, \dots$
 - compute π_{v_k} and P_k
 - $\mathbf{d}_k = P_k^\top \mathbf{d}_{k-1}$
 - $v_{k+1} = \mathbf{T}^{\text{R1}}(v_k) = \mathbf{T}(v_k) + \frac{\gamma}{1-\gamma} \left\langle \mathbf{d}_k^\top, \mathbf{T}(v_k) - v_k \right\rangle \mathbf{1}$
-

Power Iteration

$$\mathbf{d}_k^{(i+1)} = P_k \mathbf{d}_k^{(i)} \quad \text{for } i = 0, 1, \dots, N$$

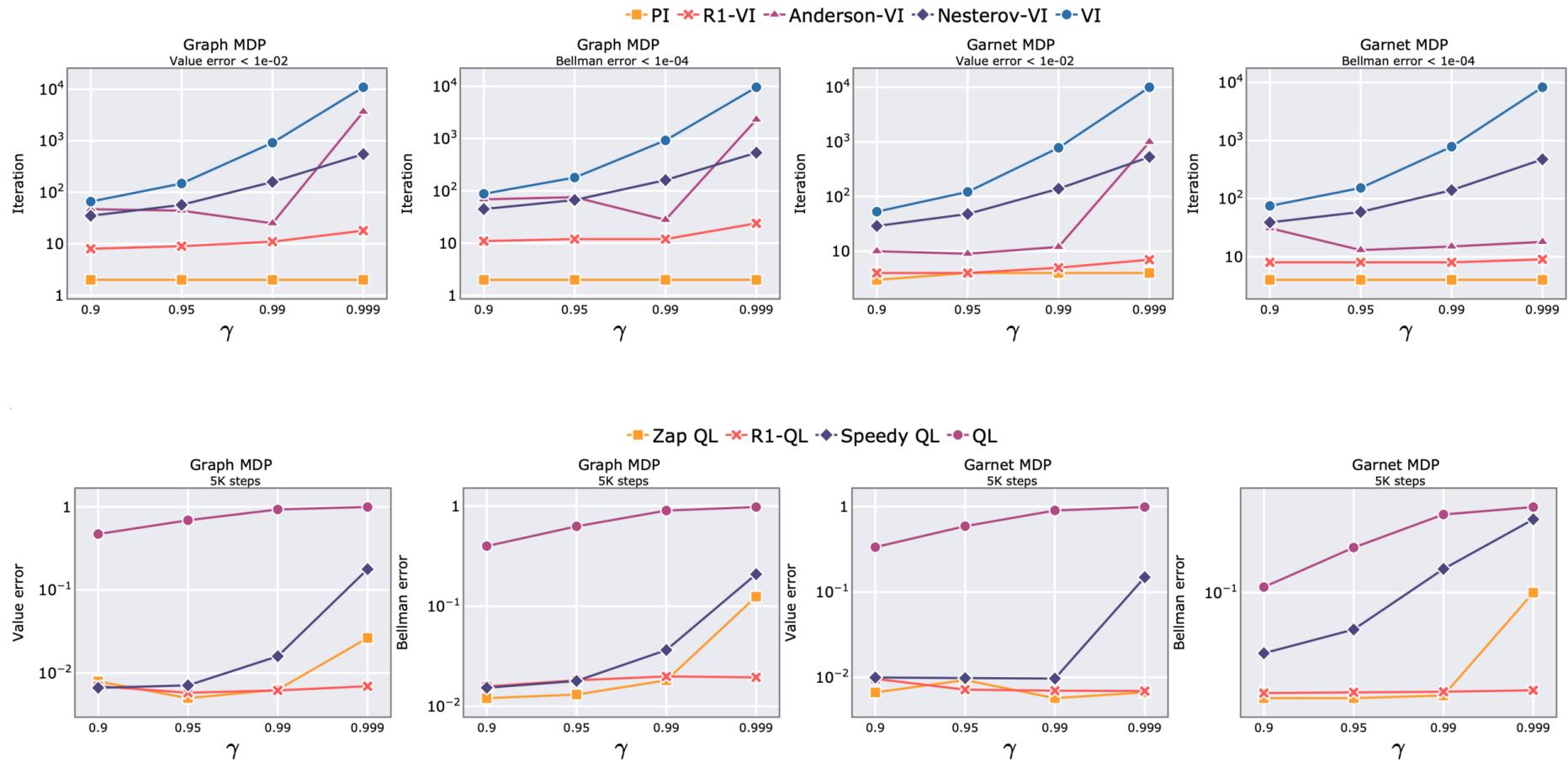
set $N \leftarrow 1$

Experiments



→ **R1-VI** requires a similar number of iterations for larger γ values.

Experiments



Summary & Limitations

- R1-VI and R1-QL has the **same** computational complexity as VI and QL, respectively.
- Both R1-VI and R1-QL are **convergent** with rate γ , but the corresponding operators are **not contractions** (Theorems 3.3 & 4.2).
- Improvement in the value space, but **no improvement** in the policy space.
- Performance depends on the spectral gap of P_k .
- Adaptable to *asynchronous* learning (**no convergence guarantee**) and average reward setting.