

# A NEW POLICY BASED RL ALGORITHM WITH REDUCED BIAS AND VARIANCE

Jayakumar Subramanian & Aditya Mahajan

ECE & CIM, McGill University and GERAD



## Introduction

Policy iteration based methods for RL may be broadly classified as:

- Monte Carlo (MC) methods [low bias, high variance]
- Temporal Difference (TD) methods [high bias, low variance]

### Our Contribution

- Monte Carlo (RMC) – a method based on MC – for infinite horizon models with a designated start state [bias comparable to MC and variance comparable to TD].
- RMC retains the advantages of the Monte Carlo approach including low bias, simplicity and ease of implementation while circumventing its key drawbacks of high variance and delayed (end of episode) updates.

## Model

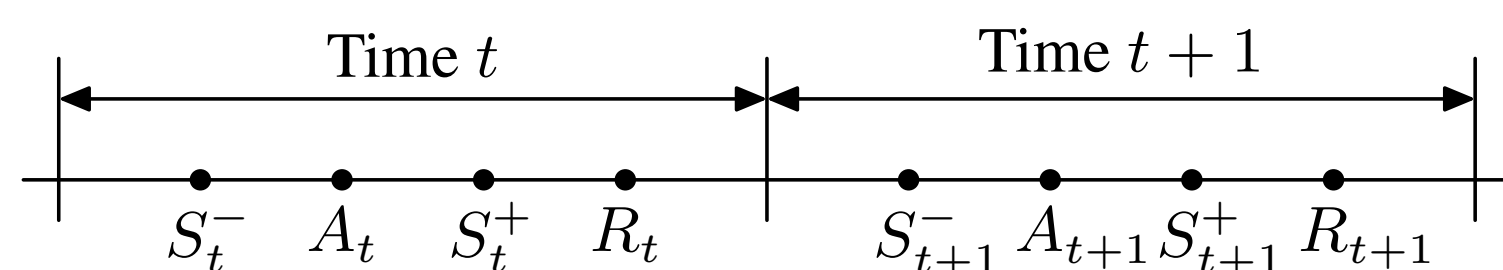


Fig. 1: A timing diagram showing the system variables.

We consider a generalized Markov decision process (MDP)  $(\mathcal{S}^-, \mathcal{S}^+, \mathcal{A}, P^-, P^+, r)$  where

- $\mathcal{S}^-$  is the pre-decision state space;
- $\mathcal{S}^+$  is the post-decision state space;
- $\mathcal{A}$  is the action space;
- $P^-$  is a Markov transition kernel from  $\mathcal{S}^+$  to  $\mathcal{S}^-$ ;
- $P^+$  is a controlled Markov transition kernel from  $\mathcal{S}^- \times \mathcal{A}$  to  $\mathcal{S}^+$ ;
- $r: \mathcal{S}^- \times \mathcal{A} \times \mathcal{S}^+ \rightarrow \mathbb{R}$  is the per-step reward function.

## Renewal Theory based Terms

- **Stopping times:**  $\tau_1 = \inf\{t > 1 : s_t^+ = s_*^+\}$ ,  $\tau_2 = \inf\{t > \tau_1 : s_t^+ = s_*^+\}$ ,
- **Discounted reward & time:**  $R_n = \gamma^{-\tau_{n-1}} \sum_{t=\tau_{n-1}+1}^{\tau_n} \gamma^{t-1} R_t$  and  $T_n = \gamma^{-\tau_{n-1}} \sum_{t=\tau_{n-1}+1}^{\tau_n} \gamma^{t-1}$ .
- By the strong Markov property,  $\{R_n\}_{n \geq 1}$  and  $\{T_n\}_{n \geq 1}$  are i.i.d. sequences. Define:  $\hat{R}_N = \frac{1}{N} \sum_{n=1}^N R_n$  and  $\hat{T}_N = \frac{1}{N} \sum_{n=1}^N T_n$ .

### Main Result

**Renewal relationship:** The post-decision value function at the reference state  $s_*^+$  is given by:

$$V_{\pi}^+(s_*^+) = \frac{R_{\pi}}{(1-\gamma)T_{\pi}}.$$

## RMC Algorithm

- Policies parameterized by a closed and convex subset  $\Theta$  of the Euclidean space.
- Given a  $\theta \in \Theta$ ,  $\pi_{\theta}$  denotes the policy parametrized by  $\theta$  and  $J_{\theta}$  to denote  $V_{\pi_{\theta}}^+(s_*^+)$ .
- Policy Improvement:  $\theta_{m+1} = [\theta_m + \alpha_m \hat{J}_{\theta_m}]_{\Theta}$
- Using renewal relationship:  $\nabla_{\theta} J_{\theta} = H_{\theta} / T_{\theta}^2$ , where  $H_{\theta} = T_{\theta} \nabla_{\theta} R_{\theta} - R_{\theta} \nabla_{\theta} T_{\theta}$ .

### Modified Policy Improvement

$$\theta_{m+1} = [\theta_m + \alpha_m \hat{H}_m]_{\Theta}; \quad \hat{H}_m \text{ is estimated as } \hat{H}_N = \hat{T}_N \hat{R}'_N - \hat{R}_N \hat{T}'_N$$

### Policy Differentiable w.r.t $\theta$

- Likelihood Ratio:  $L_t = \nabla_{\theta} \log[\pi_{\theta}(A_t | S_t^-)]$ ,  $L_n = \sum_{t=\tau_{n-1}+1}^{\tau_n} L_t$
- Estimators for  $\nabla_{\theta} R_{\theta}$  and  $\nabla_{\theta} T_{\theta}$ :  $\hat{R}'_N = \frac{1}{N} \sum_{n=1}^N R_n L_n$  and  $\hat{T}'_N = \frac{1}{N} \sum_{n=1}^N T_n L_n$ .

**Algorithm 1:** RMC Algorithm for RL: Likelihood ratio (or score function) based variant

**input** :  $\theta_0$  : Initial policy  
 $\gamma$  : Discount factor  
 $s_*^+$  : Reference state (same as initial state)  
 $N$  : Number of regenerative cycles  
 $M$  : Number of iterations

**output** :  $\theta_M$ : Estimate for best policy

**initialize:**  $t = 1$

**for** iteration  $m = 0$  **up to**  $M - 1$  **do**  
 Set  $\bar{R} = 0$ ;  $\bar{R}' = 0$ ;  $\bar{T} = 0$ ;  $\bar{T}' = 0$   
**for** regenerative cycle  $n = 1$  **up to**  $N$  **do**  
 Set  $\hat{R}_n = 0$ ;  $\hat{T}_n = 0$ ;  $L_n = 0$ ;  $\text{scale} = 1$   
**do**  
 Observe  $(S_t^-, A_t, S_t^+, R_t)$  where  $A_t \sim \theta_m(S_t^-)$   
 $\hat{R}_n \leftarrow \hat{R}_n + \text{scale} * R_t$ ;  $\hat{T}_n \leftarrow \hat{T}_n + \text{scale}$ ;  $L_n \leftarrow L_n + \nabla_{\theta} \log[\pi_{\theta}(A_t | S_t^-)]$   
 $\text{scale} \leftarrow \text{scale} * \gamma$ ;  
 $t \leftarrow t + 1$   
**while**  $S_t^+ \neq s_*^+$   
 $\bar{R} \leftarrow \bar{R} + \frac{1}{n}(\hat{R}_n - \bar{R})$ ;  $\bar{R}' \leftarrow \bar{R}' + \frac{1}{n}(\hat{R}_n L_n - \bar{R}')$   
 $\bar{T} \leftarrow \bar{T} + \frac{1}{n}(\hat{T}_n - \bar{T})$ ;  $\bar{T}' \leftarrow \bar{T}' + \frac{1}{n}(\hat{T}_n L_n - \bar{T}')$   
 $\hat{H}'_m = \bar{T} \bar{R}' - \bar{R} \bar{T}'$   
 $\theta_{m+1} = [\theta_m + \alpha_m \hat{H}'_m]_{\Theta}$   
**return**  $\theta_M$

### Policy Not Differentiable w.r.t $\theta$ – Simultaneous Perturbation FD Methods

- SPSA:  $\hat{R}'_{\theta} = (R_{\theta+\beta\delta} - R_{\theta})/\beta\delta$ , where  $\delta$ : a  $\theta$ -dim Bernoulli/Binomial random variable.
- SF:  $\hat{R}'_{\theta} = \eta(R_{\theta+\beta\eta} - R_{\theta})/\beta$ , where  $\eta$ : a  $\theta$ -dim Gaussian ( $\mathcal{N}(0, 1)$ ) random variable.
- $\beta$  is a small positive value and expressions for  $\hat{T}'_N$  are similar.

## Numerical Examples

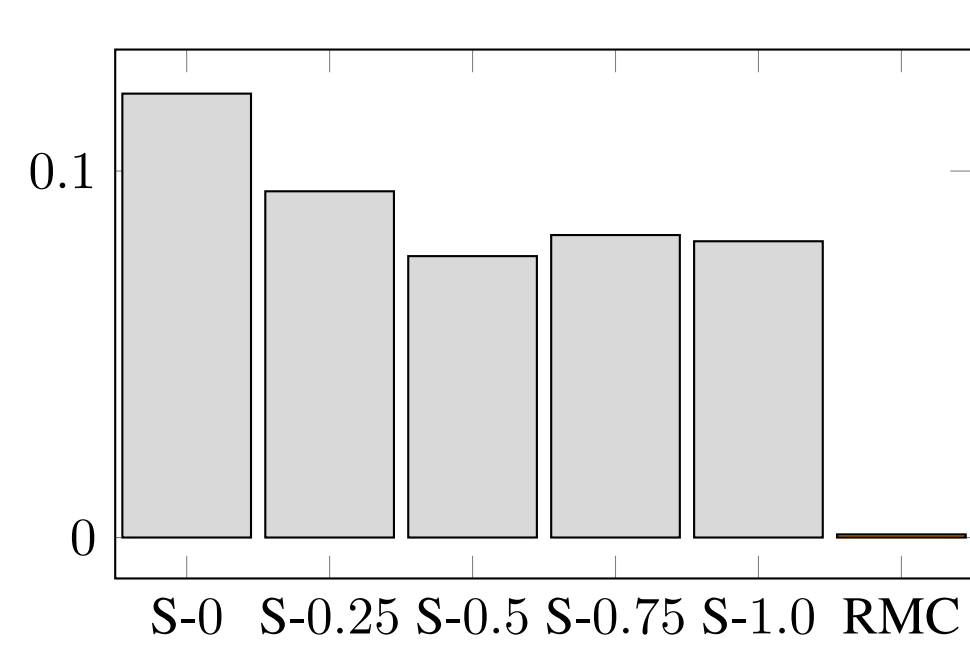


Fig. 2: Bias for Howard Taxi

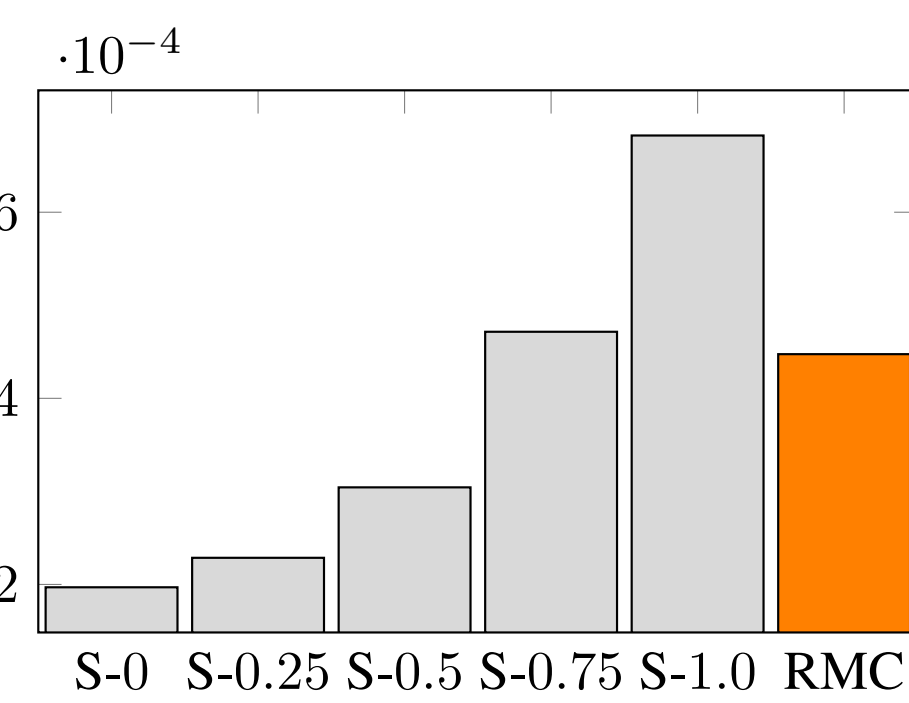


Fig. 3: Variance for Howard Taxi

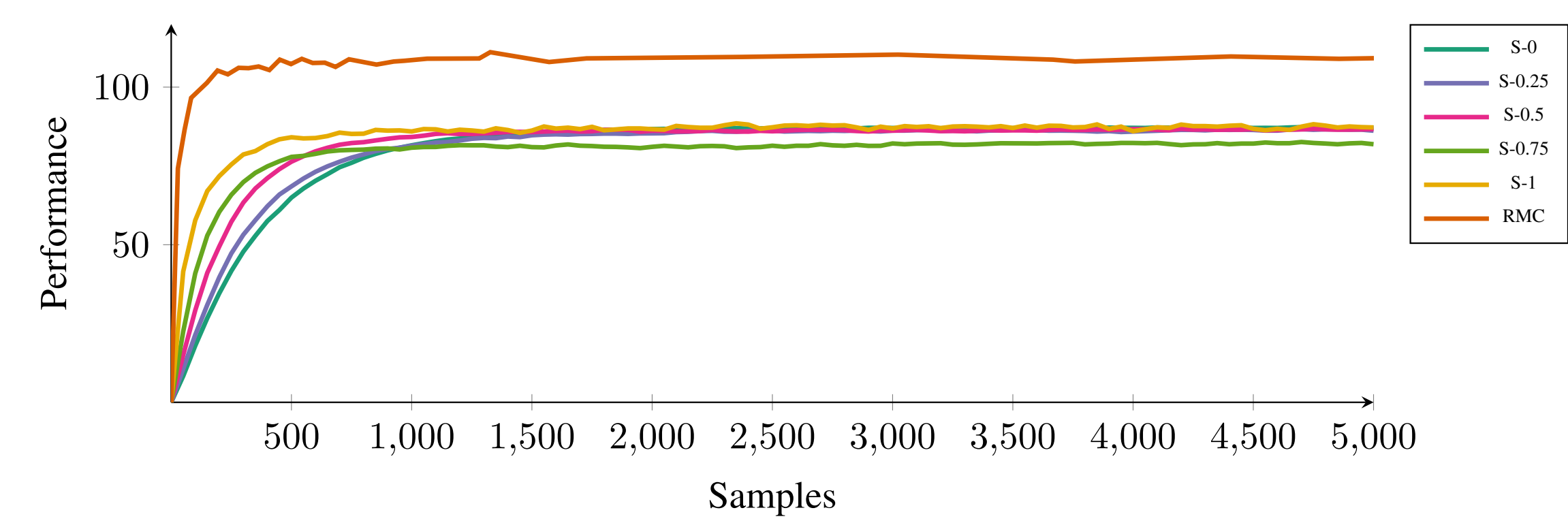


Fig. 4: Comparison of policy improvement for Howard Taxi using SARSA- $\lambda$

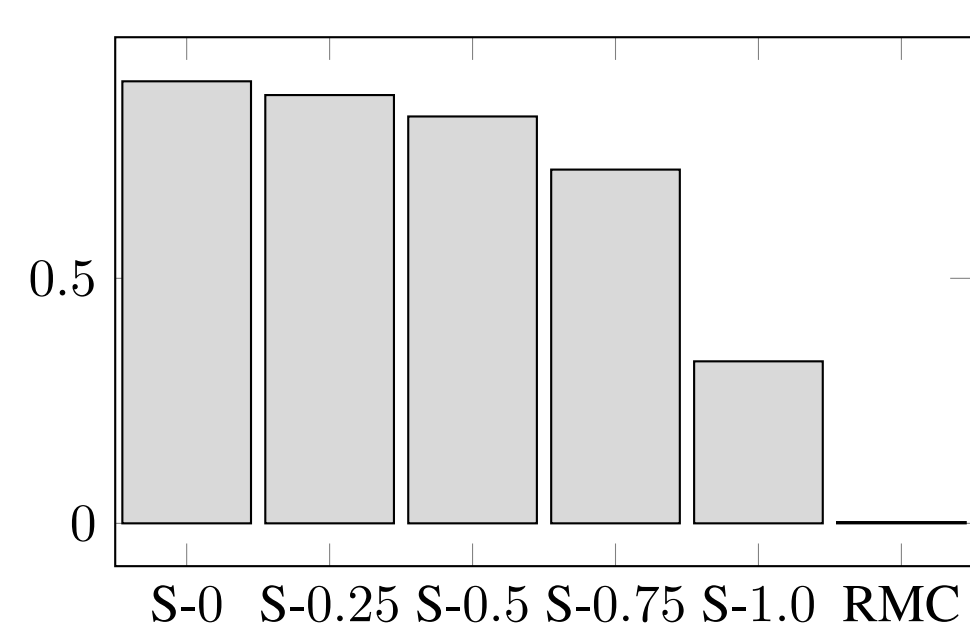


Fig. 5: Bias for Random MDP

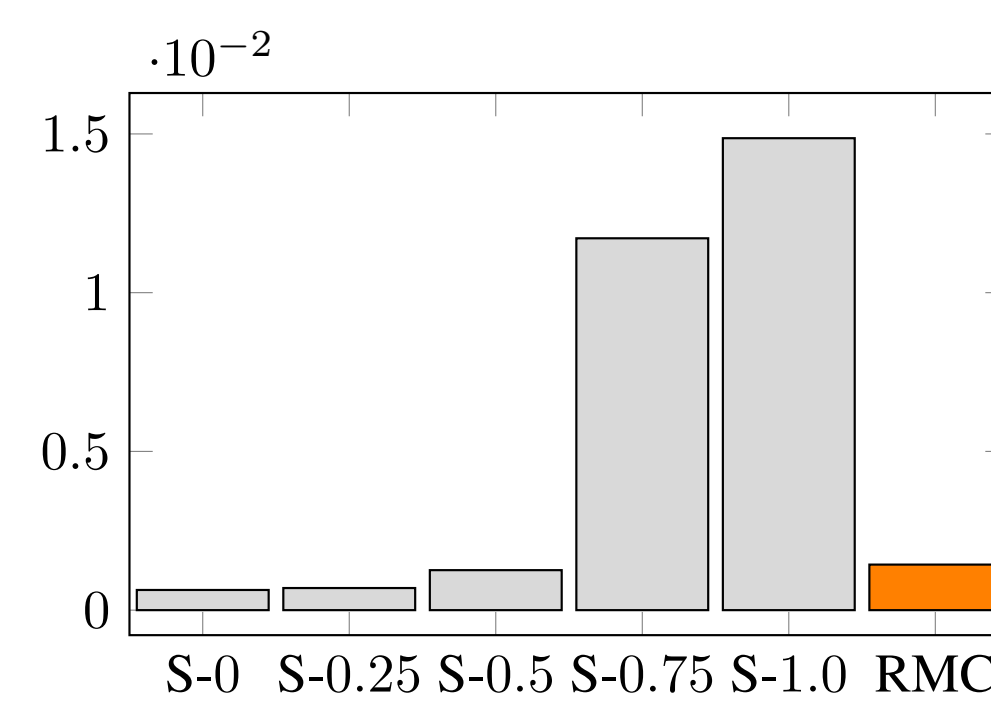


Fig. 6: Variance for Random MDP

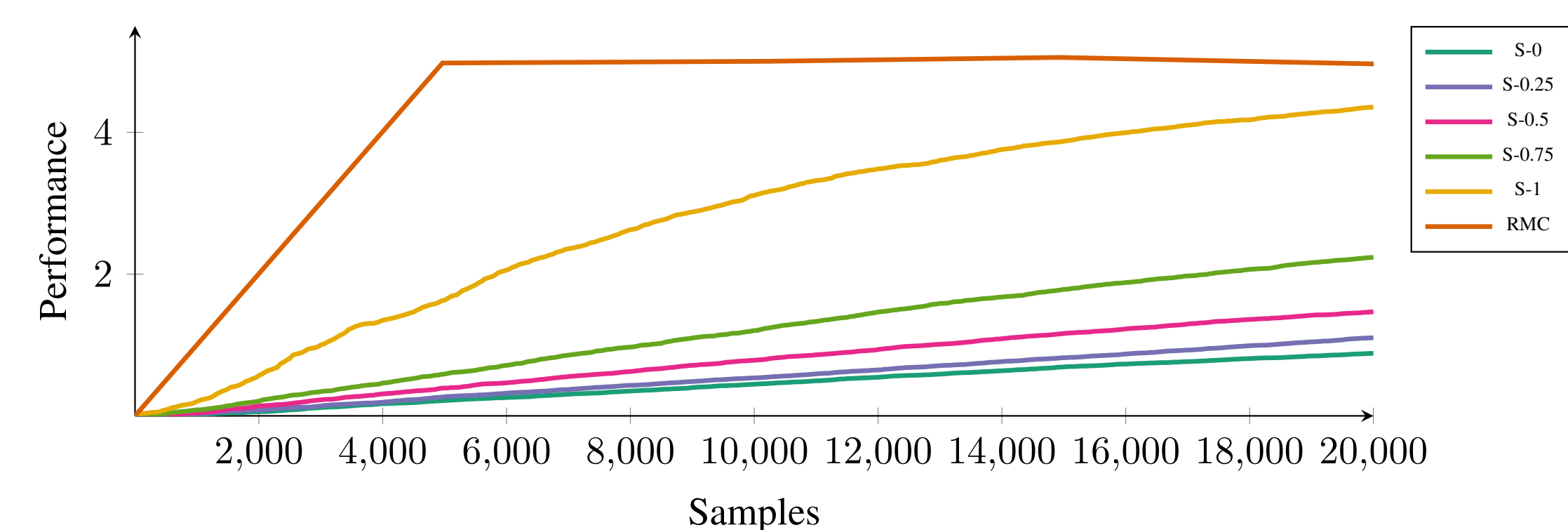


Fig. 7: Comparison of policy improvement for Random MDP using SARSA- $\lambda$

## Conclusions

- A new RL algorithm, called Renewal Monte Carlo (RMC), for infinite-horizon discounted reward problems with a designated state state
- Experimental study suggests that RMC has significantly low bias with variance that is similar to SARSA- $\lambda$
- RMC is an Actor only method and works for models with continuous state and action spaces without the need to approximate the value function and for average rewards
- Possible to obtain an “every step” variant of RMC that can be used to to estimate the entire value function (or its approximation)