Model-Free Approach to Evaluate Reinforcement Learning Algorithms

Denis Belomestny¹² Ilya Levin² Eric Moulines³² Alexey Naumov² Sergey Samsonov² Veronika Zorina²

1. Introduction

The key objective of Reinforcement Learning (RL) is to learn an optimal agent's behaviour in an unknown environment. A natural performance metric is given by the value function V^{π} which is the expected total reward of the agent following π . Unfortunately, even a precise knowledge of V^{π} does not provide information on how far is the policy π from the optimal one. To address this issue a popular quality measures are the regret bounds of the algorithm (Jaksch et al., 2010) and suboptimality gap (policy error) (Szepesvári, 2010; Pires & Szepesvári, 2016). However, available estimates of both quantities are typically pessimistic and rely on the unknown quantities of the underlying Markov Decision Processes (MDP). Moreover, even if the bounds are known, they does not apply to the general policy π and depends significantly on the particular algorithm which produced it (Jin et al., 2018; Azar et al., 2017).

In this paper we are interested in deriving agnostic (model independent) bounds for the policy error using the concept of *upper solutions* to the Bellman optimality equation. Our approach is substantially different from the ones known in literature as it can be used to estimate the suboptimality gap for an arbitrary given policy π . The concept of upper solutions is closely related to martingale duality in optimal control and information relaxation approach, see (Belomestny & Schoenmakers, 2018), and references therein. The concept of upper solutions has also a connection to distributional RL, as it can be formulated pathwise or using distributional Bellman operator, see e.g. (Lyle et al., 2019).

The contributions of this paper are three-fold:

- We propose a novel approach to construct model free confidence bounds for the optimal value function V^{*} based on a notion of upper solutions.
- Given a policy π, we propose an upper value iterative procedure (UVIP) for constructing an a.s. upper bound for V^π such that it coincides with V^{*} if π = π^{*}.

¹Duisburg-Essen University, Germany ²National Research University Higher School of Economics, Russia ³Ecole Polytechnique, France. Correspondence to: Denis Belomestny <denis.belomestny@uni-due.de>, Veronika Zorina <vozorina@hse.ru>.

Proceedings of the 38th International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

• We study convergence properties of the proposed algorithm. In particular, we show that the variance of the resulting upper bound is small if π is close to π^* leading to the tight confidence bounds for V^* .

Definitions and notations A Markov Decision Process (MDP) is a tuple (X, A, \mathcal{P}, r, γ), where X is the state space, A is the action space, $\mathscr{P} = (\mathsf{P}^a)_{a \in \mathsf{A}}$ is the *transition proba*bility kernel, $r = (r^a)_{a \in A}$ is the reward function, $0 < \gamma < 1$ is the discount factor. A policy, is denoted as π . An optimal policy π^* is one that achieves the maximum possible value amongst all policies in each state $x \in X$. The *optimal value* for state x is denoted by $V^*(x)$. The value function of a policy π in a state $x \in X$ is denoted by $V^{\pi}(x)$, similarly in a state $x \in X$ and $a \in A$ the *action-value* function $Q^{\pi}(x, a)$. Let us denote the space of bounded measurable functions with domain X by $\mathcal{B}(X)$ equipped with the norm $||f||_{\mathsf{X}} = \sup_{x \in \mathsf{X}} |f(x)|$ for any $f \in \mathcal{B}(\mathsf{X})$. The *Bellman return operator* w.r.t P, $T_P : \mathcal{B}(X) \to \mathcal{B}(X \times A)$, is defined by $(\mathsf{T}_{\mathsf{P}}V)(x,a) = r^{a}(x) + \gamma \mathsf{P}^{a}V(x)$, where $(\mathsf{P}^a V)(x) = \int V(x') \mathsf{P}^a(dx'|x)$. We also define the max*imum selection* operator M : $\mathcal{B}(X \times A) \rightarrow \mathcal{B}(X)$ by $(\mathsf{M}V^{\cdot})(x) = \max_{a} V^{a}(x)$. Then MT_{P} corresponds to the Bellman optimality operator. The optimal value function V^{\star} satisfies a non-linear fixed-point equation

$$V^{\star}(x) = \mathsf{MT}_{\mathsf{P}}V^{\star}(x). \tag{1}$$

which is the *Bellman optimality equation*. We write $Y^{x,a}, x \in X, a \in A$ for a random variable generated according to $\mathsf{P}^a(\cdot|x)$, and define a random Bellman operator $(\widetilde{\mathsf{T}}_{\mathsf{P}}V)(x) \mapsto r^a(x) + \gamma V(Y^{x,a})$.

2. UVIP Algorithm

A straightforward approach to bound the policy error $\Delta_{\pi}(x) \stackrel{\text{def}}{=} V^{*}(x) - V^{\pi}(x)$ requires the estimation of the optimal value function $V^{*}(x)$. Unfortunately, (1) does not allow to represent V^{*} as an expectation. Thus the problem of estimating V^{*} can not be naturally reduced to a stochastic approximation problem. Moreover, for sequence of value iteration procedure $V_{k+1} = \mathsf{MT}_{\mathsf{P}}V_k$, if $(\mathsf{P}^a)_{a\in\mathsf{A}}$ is replaced by its empirical estimate $\hat{\mathsf{P}}^a$ the desired upper biasedness property $V_k(x) \geq V^{*}(x)$ is lost. Below we describe our approach, which is based on the following key assumptions:

- we consider infinite-horizon MDPs with discount factor γ < 1;
- we can sample from the conditional distribution $P^{a}(\cdot|x)$ for any $x \in X$ and $a \in A$.

The key concept of our algorithm is an *upper solution*, introduced below.

Definition 2.1. We call a function V^{up} an upper solution to the Bellman optimality equation (1) if

$$V^{\mathrm{up}}(x) \ge \mathsf{MT}_{\mathsf{P}}V^{\mathrm{up}}(x), \forall x \in \mathsf{X}.$$

Upper solutions can be used to build tight upper bounds for the optimal value function V^* . Let $\Phi \in \mathcal{B}(X)$ be a *martingale function* w.r.t. the operator P^a , that is, $\mathsf{P}^a\Phi(x) = 0$ for all $a \in \mathsf{A}, x \in \mathsf{X}$. Define V^{up} as a solution to the following fixed point equation:

$$V^{\text{up}}(x) = \mathbb{E}[\max_{a} \{r^{a}(x) + \gamma(V^{\text{up}}(Y^{x,a}) - \Phi(Y^{x,a}))\}], \quad (2)$$

where $Y^{x,a} \sim \mathsf{P}^{a}(\cdot|x)$. In terms of the random Bellman operator $\widetilde{\mathsf{T}}_{\mathsf{P}}$, we can rewrite (2) as $V^{\mathrm{up}} = \mathbb{E}[\mathsf{M}\widetilde{\mathsf{T}}_{\mathsf{P}}(V^{\mathrm{up}} - \Phi)]$. It is easy to see that (2) defines an upper solution. Indeed, for any $x \in \mathsf{X}$,

$$V^{\mathrm{up}}(x) \geq \max_{a} \mathbb{E}[r^{a}(x) + \gamma(V^{\mathrm{up}}(Y^{x,a}) - \Phi(Y^{x,a}))]$$

$$= \max_{a} \{r^{a}(x) + \gamma \mathsf{P}^{a} V^{\mathrm{up}}(x)\} = \mathsf{MT}_{\mathsf{P}} V^{\mathrm{up}}(x)$$

Note that unlike the optimal state value function V^* , the upper solution V^{up} is represented as an expectation, which allows us to use various stochastic approximation methods to compute V^{up} . The Banach's fixed-point theorem implies that for iterates

$$V_{k+1}^{\rm up} = \mathbb{E}[\mathsf{M}\widetilde{\mathsf{T}}_{\mathsf{P}}(V_k^{\rm up} - \Phi)], \quad k \in \mathbb{N},$$

we have convergence $V_k^{\text{up}} \to V^{\text{up}}$ as $k \to \infty$. Moreover, V^{up} does not depend on V_0^{up} and $V_k^{\text{up}}(x) \ge V^*(x)$ for any $k \in \mathbb{N}, x \in X$, provided that $V_0^{\text{up}}(x) \ge V^*(x)$. Given a policy π and the corresponding value function V^{π} , we set $\Phi_{\pi}^{x,a}(y) \stackrel{\text{def}}{=} V^{\pi}(y) - (\mathbb{P}^a V^{\pi})(x)$. It is easy to check that $\mathbb{P}^a \Phi_{\pi}^{x,a}(x) = 0$. This leads to the upper value iterative procedure (UVIP):

$$V_{k+1}^{\rm up}(x) = \mathbb{E}[\mathsf{M}\widetilde{\mathsf{T}}_{\mathsf{P}}(V_k^{\rm up} - \Phi_{\pi}^{x,\cdot})(x)] = \\ \mathbb{E}\Big[\max_a \{r^a(x) + \gamma(V_k^{\rm up}(Y^{x,a}) - \Phi_{\pi}^{x,a}(Y^{x,a}))\}\Big]$$
(3)

with $V_0^{\text{up}} \in \mathcal{B}(\mathsf{X})$. Further note that by taking $\Phi^{x,a}(y) \stackrel{\text{def}}{=} V^*(y) - (\mathsf{P}^a V^*)(x)$, we get with probability 1 :

$$V^{\star}(x) = (\mathsf{M}\widetilde{\mathsf{T}}_{\mathsf{P}}(V^{\star} - \Phi^{x,\cdot}))(x) = \max_{a} \{r^{a}(x) + \gamma(V^{\star}(Y^{x,a}) - \Phi^{x,a}(Y^{x,a}))\}, \quad (4)$$

that is, (4) can be viewed as an almost sure version of the Bellman equation $V^* = MT_PV^*$. The upper solutions can be used to evaluate the quality of the policies and to construct confidence intervals for V^* . It is clear that

$$V^{\pi}(x) \le V^{\star}(x) \le V_k^{\mathrm{up}}(x)$$

for any $k \in \mathbb{N}$ and $x \in \mathsf{X}$, thus a policy π can be evaluated by computing the difference $\Delta_{\pi,k}^{\mathrm{up}}(x) \doteq V_k^{\mathrm{up}}(x) - V^{\pi}(x) \ge \Delta_{\pi}(x)$. Representations (3) and (4) imply

$$\|V_{k+1}^{up} - V^{\star}\|_{\mathsf{X}} \le \gamma \|V_{k}^{up} - V^{\star}\|_{\mathsf{X}} + 2\gamma \|V^{\pi} - V^{\star}\|_{\mathsf{X}},$$

 $k \in \mathbb{N}$. Hence, we derive that $\Delta_{\pi}^{up} \doteq \lim_{k \to \infty} \Delta_{\pi,k}^{up}$ satisfies

$$\|\Delta_{\pi}\|_{\mathsf{X}} \le \|\Delta_{\pi}^{\mathrm{up}}\|_{\mathsf{X}} \le \left(1 + 2\gamma(1-\gamma)^{-1}\right)\|V^{\star} - V^{\pi}\|_{\mathsf{X}}.$$
 (5)

As a result $\Delta_{\pi}^{up} = 0$ if $\pi = \pi^*$ and the corresponding confidence intervals collapses into one point. The quantity $\Delta_{\pi,k}^{up}$ can be used to measure the quality of policies π obtained by many well-known algorithms like Reinforce, A2C, etc.

For simplicity, below we will describe all the results for finite state and action spaces $(|X|, |A| \le \infty)$, providing a short remark on a generalization of these results to continuous ones. Basically, the general iteration procedure is given by (3).

For all expectations in (3) we use empirical counterparts. Algorithm 1 contains the pseudocode of UVIP.

Algorithm 1 UVIP **Input:** $V^{\pi}, \widehat{V}_0^{up}, \gamma, \varepsilon, M_1, M_2$ Output: V^{up} for $x \in \mathsf{X}, a \in \mathsf{A}$ do $\overline{V}(x, a) = M_1^{-1} \sum_{i=1}^{M_1} V^{\pi}(Y_i^{x, a}), \ Y_i^{x, a} \sim \mathsf{P}^a(\cdot|x)$ for $y \in X$ do $\Phi_{\pi}^{x,a}(y) = V^{\pi}(y) - \overline{V}(x,a)$ end for end for k = 1while $\|\widehat{V}_k^{\mathrm{up}} - \widehat{V}_{k-1}^{\mathrm{up}}\|_{\mathsf{X}} > \varepsilon$ do for $x \in \mathsf{X}$ do $$\begin{split} \widehat{V}^{\rm up}_{k+1}(x) &= M_2^{-1} \sum_{i=1}^{M_2} [\max_a \{r^a(x) + \\ &+ \gamma(\widehat{V}^{\rm up}_k(Y^{x,a}_i) - \Phi^{x,a}_{\pi}(Y^{x,a}_i))\}], \\ &Y^{x,a}_i \sim \mathsf{P}^a(\cdot|x) \end{split}$$ end for k = k + 1end while $V^{\mathrm{up}} = \widehat{V}_k^{\mathrm{up}}$

3. Convergence results

In this section, we analyze the distance between $(\widehat{V}_k^{up})_{k \in \mathbb{N}}$ and V^* , where $\widehat{V}_k^{up}(x)$ is the *k*-th iterate of Algorithm 1. Note that with $\overline{V}_{k}^{\text{up}}(x) \stackrel{\text{def}}{=} \mathsf{E}[\widehat{V}_{k}^{\text{up}}(x)]$ we have

$$\overline{V}_{k}^{\mathrm{up}}(x) \ge \max_{a} \big\{ r^{a}(x) + \gamma \mathsf{P}^{a} \overline{V}_{k-1}^{\mathrm{up}}(x) \big\}, \qquad (6)$$

 $x \in X, \quad k \in \mathbb{N}.$ Furthermore, if $\widehat{V}_0^{\mathrm{up}}(x) \geq V^{\star}(x)$ for $x \in X$, then $\overline{V}_k^{\mathrm{up}}(x) \geq V^{\star}(x)$ for any $x \in X$ and $k \in \mathbb{N}$. Hence $\widehat{V}_k^{\mathrm{up}}$ is an upper-biased estimate of V^{\star} for any $k \geq 0$. Before stating our convergence results, we first state a number of technical assumptions.

A1. There exists a measurable mapping $\psi : X \times A \times \mathbb{R}^m \to X$ such that $Y^{x,a} = \psi(x, a, \xi)$, where ξ is a random variable with values in $\Xi \subseteq \mathbb{R}^m$ and distribution P_{ξ} on Ξ , that is, $\psi(x, a, \xi) \sim \mathsf{P}^a(\cdot|x)$.

A2. For some $R_{\max} > 0$ and all $a \in A$, $||r^a||_X \le R_{\max}$.

Suppose that for each k = 1, ..., K we use an i.i.d. sample $\boldsymbol{\xi}_k = (\xi_{k,1}, ..., \xi_{k,M_1+M_2}) \sim \mathbf{P}_{\boldsymbol{\xi}}^{\otimes (M_1+M_2)}$ to generate $Y_j^{x,a} = \psi(x, a, \xi_{k,j}), j = 1, ..., M_1 + M_2$, and these samples are independent for different k. We now state main theorems that can be proved.

Theorem 3.1. Assume A1, A2. Then for any $k \in \mathbb{N}$ and $\delta \in (0, 1)$ it holds with probability at least $1 - \delta$ that

$$\|\widehat{V}_{k}^{\text{up}} - V^{*}\|_{\mathsf{X}} \lesssim \gamma^{k} \|\widehat{V}_{0}^{\text{up}} - V^{*}\|_{\mathsf{X}} + \|V^{\pi} - V^{*}\|_{\mathsf{X}} + \sqrt{\frac{\log(|\mathsf{X}||\mathsf{A}|/\delta)}{M_{1}}} \,.$$
(7)

In the above bound \lesssim stands for inequality up to a constant depending on γ and R_{max} .

Variance of the estimator and confidence bounds. Our next step is to bound the variance of the estimator $\hat{V}_k^{\mathrm{up}}(x)$. Denote

$$\sigma_k \stackrel{\text{def}}{=} \gamma^k \| \widehat{V}_0^{\text{up}} - V^* \|_{\mathsf{X}} + \| V^\pi - V^* \|_{\mathsf{X}} + \sqrt{\frac{\log(|\mathsf{X}||\mathsf{A}|)}{M_1}} \,. \tag{8}$$

Note that under A1, A2, and $\left\| \widehat{V}_0^{\text{up}} \right\|_{\mathsf{X}} \leq R_{\max}(1-\gamma)^{-1}$,

$$\sigma_k \lesssim \|V^{\pi} - V^{\star}\|_{\mathsf{X}} , \qquad (9)$$

provided that k and M_1 are large enough. The next theorem implies that $\operatorname{Var}\left[\widehat{V}_k^{\operatorname{up}}(x)\right]$ can be much smaller than the standard rate $1/M_2$, provided that V^{π} is close to V^* and M_1, K are large enough.

Theorem 3.2. Assume A1, A2. Then for any $k \in \mathbb{N}$,

$$\max_{x \in \mathsf{X}} \mathsf{Var}\big[\widehat{V}_k^{\mathrm{up}}(x)\big] \lesssim \sigma_k^2 \log(\mathbf{e} \vee \sigma_\mathbf{k}^{-1}) \mathbf{M}_2^{-1} \,, \qquad (10)$$

Corollary 3.1. Recall that $\widehat{V}_k^{\text{up}}$ is an upper biased estimate of V^* in a sense that $\overline{V}_k^{\text{up}}(x) \ge V^*(x)$ provided $\widehat{V}_0^{\text{up}}(x) \ge V^*(x)$ for $x \in X$. Together with Theorem 3.2, it implies that for any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$V^{\pi}(x) \leq V^{\star}(x) \leq \widehat{V}_{k}^{\mathrm{up}}(x) + \sigma_{k} \sqrt{\mathrm{C}\log(\mathrm{e} \vee \sigma_{k}^{-1})\delta^{-1}M_{2}^{-1}}.$$
 (11)

Note that bounds of type (11) are known in the literature only in the case of specific policies π . For example, (Wainwright, 2019) proves bounds of this type for greedy policies in tabular Q-learning. At the same time, (11) holds for arbitrary policy π .

4. Numerical Results

In this section we demonstrate the performance of Algorithm 1 on several tabular state-space RL problems. We construct the upper confidence bounds for policy π , coming from the Value iteration procedure and Reinforce algorithm. Recall that the closer policy π is to the optimal one π^* , the smaller is the difference between $V^{\pi}(x)$ and $V^{\text{up},\pi}(x)$.



Figure 1. The difference between $V^{up,\pi_i}(x)$ and $V^{\pi_i}(x)$. X-axis represents states in a discrete environment for all pictures. Each group of three pictures of the same color demonstrates the process of learning the policy from the first iteration to the last. *First and second rows*: Evaluation of the policies during the process of Value iteration for Garnet and Frozen Lake. The policies are obtained greedily from $Q_i(x, a)$ function at the i-th step. *Third row*: Reinforce algorithm during learning for Frozen Lake.

We consider 2 popular tabular environments: Garnet ((Archibald et al., 1995)) and AI Gym Frozen Lake ((Brockman et al., 2016)). For each environment we perform Kupdates of the Value iteration with known transition kernel P^a . We denote the *k*-th step estimate of the action-value function as $\hat{Q}_k(x, a)$ and denote π_k the greedy policy w.r.t. $\hat{Q}_k(x, a)$. Then we evaluate the policies π_k with the Algorithm 1 for certain iteration numbers *k*. Figure 1 displays the gap between $V^{\pi_k}(x)$ and $V^{\text{up},\pi_k}(x)$, which converges to zero while π^k converges to the optimal policy π^* . Data for the upper bounds estimation is generated using *off-policy* method. On the Frozen Lake environment we also apply the tabular version of the **Reinforce** algorithm. We evaluate policies π_k obtained from the k-th **Reinforce** iteration. On the Figure 1 we display $V^{\pi_k}(x)$ and $V^{\text{up},\pi_k}(x)$ for different time steps k. The difference $V^{\text{up},\pi_k}(x) - V^{\pi_k}(x)$ does not converge to zero, indicating suboptimality of the **Reinforce** policy.

5. Conclusion

We propose a new approach towards model-free evaluation of the agent's policies in RL, based on upper solutions to the Bellman optimality equation (1). To the best of our knowledge, the UVIP is the first procedure which allows to construct the non-asymptotic confidence intervals for the optimal value function V^* based on the value function corresponding to an arbitrary policy π . In our analysis we consider only infinite-horizon MDPs and assume that sampling from the conditional distribution $P^a(\cdot|x)$ is feasible for any $x \in X$ and $a \in A$. A promising future research direction is to generalize the algorithm to the case of finite-horizon MDPs combining it with the idea of Real-time dynamic programming (see (Efroni et al., 2019)).

It is worth to highlight that the Theorems 3.1 and 3.2 have a generalization for the case of infinite state and action spaces, which requires the introduction of the covering number of set $X \times A$, the Dudley's integral, along with the proper approximation for $V^{up}(x)$. Moreover, the UVIP can be adapted for RL benchmarks with continuous state space and discrete action space by performing an additional approximation step. It implies that for an arbitrary policy we can construct the upper confidence bounds for such environments as AI Gym CartPole and Acrobot. Nevertheless, the success of the procedure relies on the policy evaluation methods. We have to choose the approximation points properly to be able to assess the next states quality after one step of the agent.

References

- Archibald, T. W., McKinnon, K. I. M., and Thomas, L. C. On the generation of Markov Decision Processes. *The Journal of the Operational Research Society*, 46(3):354– 361, 1995. ISSN 01605682, 14769360.
- Azar, M. G., Osband, I., and Munos, R. Minimax regret bounds for reinforcement learning. *Proceedings of the* 34th International Conference on Machine Learning, 70: 263–272, 06–11 Aug 2017.

Belomestny, D. and Schoenmakers, J. Advanced Simulation-

Based Methods for Optimal Stopping and Control: With Applications in Finance. Springer, 2018.

- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym, 2016.
- Efroni, Y., Merlis, N., Ghavamzadeh, M., and Mannor, S. Tight regret bounds for model-based reinforcement learning with greedy policies. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Jaksch, T., Ortner, R., and Auer, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(51):1563–1600, 2010.
- Jin, C., Allen-Zhu, Z., Bubeck, S., and Jordan, M. I. Is q-learning provably efficient? In Advances in Neural Information Processing Systems, volume 31. Curran Associates, Inc., 2018.
- Lyle, C., Bellemare, M. G., and Castro, P. S. A comparative analysis of expected and distributional reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 4504–4511, 2019.
- Pires, B. Á. and Szepesvári, C. Policy error bounds for model-based reinforcement learning with factored linear models. In *Conference on Learning Theory*, pp. 121–151. PMLR, 2016.
- Szepesvári, C. Algorithms for reinforcement learning. *Synthesis lectures on artificial intelligence and machine learning*, 4(1):1–103, 2010.
- Wainwright, M. J. Variance-reduced q-learning is minimax optimal. CoRR, 2019.