

UVIP: Model-Free Approach to Evaluate Reinforcement Learning Algorithms

Denis Belomestny

Duisburg-Essen University
HSE University



NATIONAL RESEARCH
UNIVERSITY

July 16, 2021

Joint work with



Sergey Samsonov



Ilya Levin



Eric Moulines



Alexey Naumov



Veronika Zorina

Paper available at <https://arxiv.org/abs/2105.02135>

Markov Decision Process (MDP)

- ▶ X - state space. By $(X_k)_{k \geq 0}$ we denote a sequence of random states.
- ▶ A - action space. Let $(A_k)_{k \geq 0}$ be a sequence of random actions.
- ▶ Agent's policy π is the distribution on A

$$\pi(a|x) = \mathbb{P}(A_k = a | X_k = x)$$

- ▶ Family of Markov transition kernels $(P^a(x'|x))_{a \in A}$:

$$P^a(x'|x) := \mathbb{P}(X_k = x' | X_{k-1} = x, A_{k-1} = a).$$

- ▶ (Deterministic) reward $r^a(x) : A \times X \mapsto \mathbb{R}$
- ▶ At step k in the state $X_k = x$ the agent performs an action $A_k \sim \pi(\cdot|x)$ (suppose $A_k = a$), obtains a reward $r^a(X_k)$ and transits to $X_{k+1} \sim P^a(\cdot|x)$

Markov Decision Process

Let $\gamma \in (0, 1]$ be the discount factor. Tuple (X, A, P, r, γ) is called Markov Decision Process.

How to measure policy's quality?

Given policies π_1 and π_2 , how to compare their quality?

How to measure policy's quality?

Given policies π_1 and π_2 , how to compare their quality?

Value function, associated with the policy π , is defined as

$$V^\pi(x) := \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r^{A_k}(X_k) \mid X_0 = x \right]$$

How can we estimate this quantity?

1. **Monte-Carlo.** Run a series of independent simulations to compute $V^\pi(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r^{A_k}(X_k) \mid X_0 = x \right]$.
2. **Stochastic Approximation.** Use Robbins-Monro procedure to obtain *Temporal difference*-based method (e.g. *TD(0)*, [Sutton \[1988\]](#), [Tsitsiklis and Van Roy \[1997\]](#)). Preliminary estimates are available even at the early stage.

Optimal policy and optimal value function

We call policy π^* an *optimal policy*, if $V^{\pi^*}(x) \geq V^\pi(x)$ for any $x \in X$.

The following result holds (see, e.g., [Puterman \[2014\]](#)):

Theorem

When the reward function is bounded, one can always find a *deterministic* Markov policy that is optimal. Moreover, the optimal value function V^* satisfies the *Bellman optimality equation*:

$$V^*(x) = \max_{a \in A} [r^a(x) + \gamma P^a V^*(x)].$$

In the formula above $P^a V^*(x) = \int_X V^*(x') P^a(dx'|x)$.

Policy quality metrics

How far is the given policy π from the optimal one π^* ?

- ▶ Popular performance metrics is the *total regret* (Azar et al. [2017], Jin et al. [2018]) of the learning algorithm with respect to an optimal policy.
- ▶ Given a sequence of policies $(\pi_k)_{k \in \mathbb{N}}$ and a sequence of episodes starting from the states $(x_0^k)_{k \in \mathbb{N}}$, we define the regret after K iterations as

$$\sum_{k=1}^K [V^*(x_0^k) - V^{\pi_k}(x_0^k)].$$

- ▶ Regret bounds are available for many tabular MDP learning algorithms. Yet the regret analysis does not provide much information for the *policy approximation methods*.

Policy quality metrics

How far is the given policy π from the optimal one π^* ?

- ▶ Consider, for example, the family of *policy approximation methods*;
- ▶ We approximate $\pi(a|x)$ by the parametric family of policies $\pi_\theta(a|x)$, $\theta \in \Theta \subseteq \mathbb{R}^{d'}$;
- ▶ It is hard to quantify the *approximation error*

$$V^*(x) - \sup_{\theta \in \Theta} V^{\pi_\theta}(x),$$

moreover, the bounds on approximation error are typically pessimistic and relies upon unknown regularity properties of given MDP;

- ▶ Is it a generic way to estimate the *policy error*

$$\Delta_\pi(x) \doteq V^*(x) - V^\pi(x),$$

if π^* (and V^*) are unknown?

Value iteration algorithm

- ▶ First idea: construct an upper biased estimate of $V^*(x)$;
- ▶ Assume that the transition kernel $(P^a)_{a \in A}$ is known;
- ▶ The *value iteration* algorithm (Bertsekas and Shreve [1978]): starting from some $V_0(x)$, iterate

$$V_{k+1} = \max_{a \in A} [r^a(x) + \gamma P^a V_k(x)] .$$

- ▶ An important property: $V_k(x) \geq V^*(x)$ for any $x \in X$ and $k \in \mathbb{N}$, provided that $V_0(x) \geq V^*(x)$.
- ▶ Unfortunately, this property is lost if $(P^a)_{a \in A}$ is not known and could be only approximated.

Model assumptions

1. Our aim is to construct an (upper biased) estimate of $\Delta_\pi(x)$;
2. We consider infinite-horizon MDPs with discount factor $\gamma < 1$;
3. We can sample from the conditional distribution $P^a(\cdot|x)$ for any $x \in X$ and $a \in A$;
4. Value function $V^\pi(x)$ is known .

Upper solutions concept

Upper solution

We call a function V^{up} an *upper solution* to the Bellman optimality equation, if

$$V^{\text{up}}(x) \geq \max_{a \in A} \{r^a(x) + \gamma P^a V^{\text{up}}(x)\}, \forall x \in X. \quad (1)$$

If $V^{\text{up}}(x)$ satisfies (1), then using the Bellman optimality equation we can show that $V^{\text{up}}(x) \geq V^*(x)$ for any $x \in X$.

The upper solution concept

- ▶ The concept is related to the martingale duality approach in optimal control;
- ▶ Rogers [2007];
- ▶ Belomestny and Schoenmakers [2018];
- ▶ Shar and Jiang [2020] used information relaxation approach to learn (approximate) upper and lower bounds for the optimal action-value function.

How to construct an upper solution?

- ▶ In practice we can construct an upper solution as follows. Consider arbitrary *martingale function* $\Phi : X \mapsto \mathbb{R}$, such that $P^a \Phi(x) = 0$ for any $a \in A, x \in X$;
- ▶ Define V^{up} as a solution to the following fixed point equation:

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

where $Y^{x,a} \sim P^a(\cdot|x)$.

Checking upper solution properties

Indeed, we easily check that V^{up} is an upper solution:

$$\begin{aligned} V^{\text{up}}(x) &= \mathbb{E}[\max_a \{r^a(x) + \gamma(V^{\text{up}}(Y^{x,a}) - \Phi(Y^{x,a}))\}] \\ &\geq \max_a \mathbb{E}[r^a(x) + \gamma(V^{\text{up}}(Y^{x,a}) - \Phi(Y^{x,a}))] \\ &= \max_a \{r^a(x) + \gamma P^a V^{\text{up}}(x)\} \end{aligned}$$

Solving the fixed point problem

- ▶ Recall that V^{up} as a solution to the following fixed point equation:

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

- ▶ Consider the iteration process

$$V_{k+1}^{\text{up}}(x) = \mathbb{E}[\max_a \{r^a(x) + \gamma(V_k^{\text{up}}(Y^{x,a}) - \Phi(Y^{x,a}))\}]$$

- ▶ The upper biasedness property is preserved: if $V_0^{\text{up}}(x) \geq V^*(x)$ for any $x \in X$, then $V_k^{\text{up}}(x) \geq V^*(x)$ for any k and x .

UVIP

- ▶ Given a policy π to be evaluated and the corresponding value function V^π , we set

$$\Phi_\pi^{x,a}(y) \doteq V^\pi(y) - (P^a V^\pi)(x).$$

- ▶ Then we write the update of the Upper value iterative procedure (UVIP) as follows:

$$V_{k+1}^{\text{up}}(x) = \mathbb{E}[\max_a \{r^a(x) + \gamma(V_k^{\text{up}}(Y^{x,a}) - \Phi_\pi^{x,a}(Y^{x,a}))\}]$$

- ▶ We enjoy the same upper biasedness property: $V_k^{\text{up}}(x) \geq V^*(x)$ for any k and x provided that $V_0^{\text{up}}(x) \geq V^*(x)$;
- ▶ Then we can evaluate the policy π by computing the difference

$$\Delta_{\pi,k}^{\text{up}}(x) \doteq V_k^{\text{up}}(x) - V^\pi(x) \geq \Delta_\pi(x).$$

Self-normalizing property

- ▶ Consider the case $\pi = \pi^*$, then the corresponding martingale function is given by

$$\Phi^{x,a}(y) \doteq V^*(y) - (P^a V^*)(x).$$

- ▶ The corresponding fixed point equation writes as

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

- ▶ It is easy to check that $V^*(x)$ in this scenario is a solution to (2), indeed,

$$\begin{aligned} V^*(x) &= \mathbb{E}[\max_a \{r^a(x) + \gamma(V^*(Y^{x,a}) - V^*(Y^{x,a}) + (P^a V^*)(x))\}] \\ &= \max_a \{r^a(x) + \gamma P^a V^*(x)\}, \end{aligned}$$

which coincides with Bellman optimality equation.

- ▶ Thus $\Delta_{\pi,k}^{\text{up}}(x) \rightarrow 0$ as $k \rightarrow \infty$.

Approximate UVIP

- ▶ Consider the $(k + 1)$ -th UVIP iteration:

$$V_{k+1}^{\text{up}}(x) = \mathbb{E}[\max_a \{r^a(x) + \gamma(V_k^{\text{up}}(Y^{x,a}) - \Phi_{\pi}^{x,a}(Y^{x,a}))\}]$$

- ▶ We need to approximate an outer expectation and the one-step transition operator P^a .

Tabular UVIP: Monte-Carlo version

Algorithm 1: UVIP

Require: $V^\pi, \widehat{V}_0^{\text{up}}, \gamma, \varepsilon, M_1, M_2$

Ensure: V^{up}

for $x \in X, a \in A$ **do**

$$\bar{V}(x, a) = M_1^{-1} \sum_{i=1}^{M_1} V^\pi(Y_i^{x,a}), \quad Y_i^{x,a} \sim P^a(\cdot|x)$$

for $y \in X$ **do**

$$\Phi_\pi^{x,a}(y) = V^\pi(y) - \bar{V}(x, a)$$

end for

end for

$k = 1$

while $\|\widehat{V}_k^{\text{up}} - \widehat{V}_{k-1}^{\text{up}}\|_X > \varepsilon$ **do**

for $x \in X$ **do**

$$\widehat{V}_{k+1}^{\text{up}}(x) = M_2^{-1} \sum_{i=1}^{M_2} [\max_a \{r^a(x) + \gamma(\widehat{V}_k^{\text{up}}(Y_i^{x,a}) - \Phi_\pi^{x,a}(Y_i^{x,a}))\}],$$
$$Y_i^{x,a} \sim P^a(\cdot|x)$$

end for

$k = k + 1$

end while

$V^{\text{up}} = \widehat{V}_k^{\text{up}}$

UVIP in general state space

Assumption A1

We suppose that (X, ρ_X) and (A, ρ_A) are compact metric spaces. Moreover, $X \times A$ is equipped with some metric ρ , such that $\rho((x, a), (x', a)) = \rho_X(x, x')$ for any $x, x' \in X$ and $a \in A$.

In this scenario we fix a set of points $X_N = (x_1, \dots, x_N)$ and aim to evaluate $\Delta_\pi(x)$ for $x \in X_N$.

UVIP in general state space

- ▶ Recall that $(k + 1)$ -th UVIP iteration writes as

$$V_{k+1}^{\text{up}}(x) = \mathbb{E} \left[\max_a \{ r^a(x) + \gamma (V_k^{\text{up}}(Y^{x,a}) - \Phi_{\pi}^{x,a}(Y^{x,a})) \} \right].$$

In order to update $\hat{V}_{k+1}^{\text{up}}(x_i)$ the algorithm requires to calculate $\hat{V}_k^{\text{up}}(Y^{x_i,a})$ with $Y^{x_i,a} \sim P^a(\cdot|x)$;

- ▶ This means that we need an additional interpolation step to calculate $\hat{V}_k(y)$ at points $y \notin (x_1, \dots, x_N)$;
- ▶ For example, for a Lipschitz function $f \in \text{Lip}(L)$ we can use *optimal central interpolant*

$$I[f](x) = (H_f^{\text{low}}(x) + H_f^{\text{up}}(x))/2, \text{ where}$$

$$H_f^{\text{low}}(x) = \max_{\ell \in \{1, \dots, N\}} (f(x_{\ell}) - L\rho_X(x, x_{\ell})),$$

$$H_f^{\text{up}}(x) = \min_{\ell \in \{1, \dots, N\}} (f(x_{\ell}) + L\rho_X(x, x_{\ell})).$$

Approximate UVIP

Input: Sample (x_1, \dots, x_N) ; $V^\pi, \tilde{V}_0^{\text{up}}, M_1, M_2, \gamma, \varepsilon$

while $\|\tilde{V}_k^{\text{up}} - \tilde{V}_{k-1}^{\text{up}}\|_\infty > \varepsilon$ **do**

for $a \in A$ **do**

for $i \in [N]$ **do**

for $j \in [M_1 + M_2]$ **do**

$\hat{V}_k^{\text{up}}(Y_j^{x_i, a}) = I[\tilde{V}_k^{\text{up}}](Y_j^{x_i, a}), Y_j^{x_i, a} \sim P^a(\cdot | x_i)$

end

$\bar{V}^{(i, a)} = M_1^{-1} \sum_{j=1}^{M_1} V^\pi(Y_j^{x_i, a});$

end

end

for $i \in [N]$ **do**

$\tilde{V}_{k+1}^{\text{up}}(x_i) =$

$M_2^{-1} \sum_{j=M_1+1}^{M_1+M_2} \max_{a \in A} \{r^a(x_i) + \gamma(\hat{V}_k^{\text{up}}(Y_j^{x_i, a}) - V^\pi(Y_j^{x_i, a}) + \bar{V}^{(i, a)})\};$

end

end

Theoretical assumptions

Assumption A2

There exists a measurable mapping $\psi : X \times A \times \mathbb{R}^m \rightarrow X$ 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, $\psi(x, a, \xi) \sim P^a(\cdot|x)$.

Assumption A3

For some positive constant R_{\max} and all $a \in A$, $\|r^a\|_X \leq R_{\max}$.

Assumption A4

For some positive constants $L_\psi \leq 1$, L_{\max} , L_π and all $a \in A$, $\xi \in \Xi$,

$$\text{Lip}_{\rho_X}(r^a(\cdot)) \leq L_{\max}, \quad \text{Lip}_\rho(\psi(\cdot, \cdot, \xi)) \leq L_\psi, \quad \text{Lip}_\rho((V^\pi \circ \psi)(\cdot, \cdot, \xi)) \leq L_\pi.$$

Tabular MDP's case

Checking A4

Let $|X| < \infty$ and $|A| < \infty$. Consider

$$\begin{aligned}\rho_X(x, x') &= \mathbb{1}_{\{x \neq x'\}} \\ \rho((x, a), (x', a')) &= \mathbb{1}_{\{(x, a) \neq (x', a')\}}.\end{aligned}$$

Then the assumption A4 holds with constants $L_\psi = 1$, $L_{\max} = R_{\max}$, and $L_\pi = R_{\max}/(1 - \gamma)$.

Measuring the distance between $\widehat{V}_k^{\text{up}}$ and V^*

Theorem 1a.

Let $|X|, |A| < \infty$ and assume A2, A3. Then for any $k \in \mathbb{N}$ and $\delta \in (0, 1)$ it holds with probability at least $1 - \delta$ that

$$\|\widehat{V}_k^{\text{up}} - V^*\|_X \lesssim \gamma^k \|\widehat{V}_0^{\text{up}} - V^*\|_X + \|V^\pi - V^*\|_X + \sqrt{\frac{\log(|X||A|/\delta)}{M_1}}.$$

Measuring the distance between $\widehat{V}_k^{\text{up}}$ and V^*

Theorem 1a.

Let $|X|, |A| < \infty$ and assume A2, A3. Then for any $k \in \mathbb{N}$ and $\delta \in (0, 1)$ it holds with probability at least $1 - \delta$ that

$$\|\widehat{V}_k^{\text{up}} - V^*\|_X \lesssim \gamma^k \|\widehat{V}_0^{\text{up}} - V^*\|_X + \|V^\pi - V^*\|_X + \sqrt{\frac{\log(|X||A|/\delta)}{M_1}}.$$

Theorem 1a can be generalized under A1 – A4. In this case there is an additional term in the r.h.s, depending on the *covering radius* of the set X_N w.r.t. X , that is,

$$\rho(X_N, X) = \max_{x \in X} \min_{k \in [N]} |x - X_k|.$$

Notations

In case of more general X and A we introduce some additional notations:

- ▶ Define $\mathcal{N}(X \times A, \rho, \varepsilon)$ the covering number of the set $X \times A$ w.r.t. metric ρ , that is, the smallest cardinality of an ε -net of $X \times A$ w.r.t. ρ ;
- ▶ Define D the diameter of $X \times A$, that is,

$$D = \text{diam}(X \times A) = \max_{(x,a),(x',a') \in X \times A} \rho((x, a), (x', a'));$$

- ▶ Then $\log \mathcal{N}(X \times A, \rho, \varepsilon)$ is the metric entropy of $X \times A$ and

$$I_D = \int_0^D \sqrt{\log \mathcal{N}(X \times A, \rho, u)} du$$

is the Dudley's integral.

Some examples

- ▶ Assume $|X| < \infty$, $|A| < \infty$.
- ▶ We bypass the approximation step. Consider $\rho((x, a), (x', a')) = \mathbb{1}_{\{(x,a) \neq (x',a')\}}$. Then

$$D = 1$$

$$I_D = \sqrt{\log(|X||A|)}.$$

Some examples

- ▶ Assume $X \subseteq [0, 1]^{d_X}$, $|A| < \infty$
- ▶ Let $\rho_X(x, x') = \|x - x'\|$, $\rho((x, a), (x', a')) = \|x - x'\| + \mathbb{1}_{\{a \neq a'\}}$. Then

$$D \leq \sqrt{d_X} + 1,$$

$$I_D \lesssim \sqrt{d_X \log |A|} + \sqrt{d_X \log d_X}.$$

Measuring the distance between $\widehat{V}_k^{\text{up}}$ and V^*

Theorem 1b.

Assume A1 – A4 and suppose that $\text{Lip}_{\rho_X}(\widehat{V}_0^{\text{up}}) \leq L_0$ with some constant $L_0 > 0$. Then for any $k \in \mathbb{N}$ and $\delta \in (0, 1)$ it holds with probability at least $1 - \delta$ that

$$\begin{aligned} \|\widehat{V}_k^{\text{up}} - V^*\|_X &\lesssim \gamma^k \|\widehat{V}_0^{\text{up}} - V^*\|_X + \|V^\pi - V^*\|_X + \frac{I_{\mathcal{D}} + D\sqrt{\log(1/\delta)}}{\sqrt{M_1}} \\ &\quad + \rho(X_N, X). \end{aligned}$$

In case $X = [0, 1]^d$ and $X_N = \{X_1, \dots, X_N\}$ being a set of N points, uniformly distributed over X , the following bound is available: for any $\delta \in (0, 1)$,

$$\rho(X_N, X) \lesssim \sqrt{d_X} \left(\frac{\log(1/\delta) \log N}{N} \right)^{1/d_X}$$

with probability at least $1 - \delta$.

Variance of the estimator

We need to assume in addition that $X \times A$ is a parametric class with the metric entropy satisfying the following assumption:

Assumption A5

There exist a constant $C_{X,A} > 1$ such that for any $\varepsilon \in (0, D)$,

$$\log \mathcal{N}(X \times A, \rho, \varepsilon) \leq C_{X,A} \log(1 + 1/\varepsilon).$$

Variance of the estimator

Theorem 2.

Let us introduce σ_k as an upper bound for $\mathbb{E}^{1/2}[\|\widehat{V}_k^{\text{up}} - V^*\|_{\mathcal{X}}^2]$, that is,

$$\sigma_k \doteq \gamma^k \|\widehat{V}_0^{\text{up}} - V^*\|_{\mathcal{X}} + \|V^\pi - V^*\|_{\mathcal{X}} + \frac{I_{\mathcal{D}} + D}{\sqrt{M_1}} + \rho(\mathcal{X}_N, \mathcal{X}).$$

Let A1 – A5 hold and assume additionally $\text{Lip}_{\rho_{\mathcal{X}}}(\widehat{V}_0^{\text{up}}) \leq L_0$ for some $L_0 > 0$. Then

$$\max_{x \in \mathcal{X}} \text{Var}[\widehat{V}_k^{\text{up}}(x)] \leq C \sigma_k^2 \log(e \vee \sigma_k^{-1}) M_2^{-1},$$

where the constant C depends on $C_{\mathcal{X}, A}$, γ , L_{\max} , L_ψ , L_π , L_0 and R_{\max} .

Confidence intervals for V^*

Theorem 3.

For any $x \in X_N$ and any $\delta \in (0, 1)$, with probability at least $1 - \delta$,

$$V^\pi(x) \leq V^*(x) \leq \widehat{V}_k^{\text{up}}(x) + \sigma_k \sqrt{C \log(e \vee \sigma_k^{-1}) \delta^{-1} M_2^{-1}}.$$

Moreover, for k and M_1 large enough,

$$\sigma_k \lesssim \|V^\pi - V^*\|_X.$$

Numerical results: Chain problem

- ▶ Chain is a finite MDP, where agent can move only (right and left);
- ▶ Chain has two terminal states at the ends. For transition to the terminal states agent receives 10 points and episode ends, otherwise the reward is equal to +1;
- ▶ $p \in (0, 1)$ - noise in the system, i.e. the agent's action A_k at state x is drawn from distribution

$$A_k \sim \begin{cases} \pi(\cdot|x), & \text{with probability } p; \\ \text{uniform distribution over } A, & \text{with probability } 1 - p. \end{cases}$$

Numerical results: Chain problem

- Consider the sequence of policies $(\pi_k)_{k=0}^{15}$, obtained via the value iteration procedure. We evaluate policies π_k by computing $\Delta_{\pi,k}^{\text{up}}(x)$ for some time steps k .

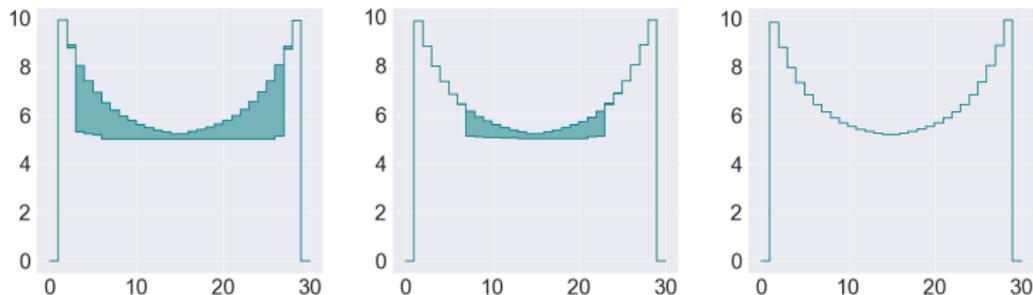


Figure: Chain environment: evaluating π_0 , π_5 and π_{15}

Numerical results: Garnet problem

- ▶ A Garnet model is specified by a triplet (N_S, N_A, N_B) , where N_S and N_A are the number of states and actions, respectively.
- ▶ The parameter N_B denotes the branching factor, that is, the number of states reachable from any state-action pair.
- ▶ We choose $N_S = 20$, $N_A = 10$, $N_B = 2$.

Numerical results: Garnet problem

- Consider the sequence of policies $(\pi_k)_{k=0}^{15}$, obtained via the value iteration procedure. We evaluate policies π_k by computing $\Delta_{\pi,k}^{\text{up}}(x)$ for some time steps k .

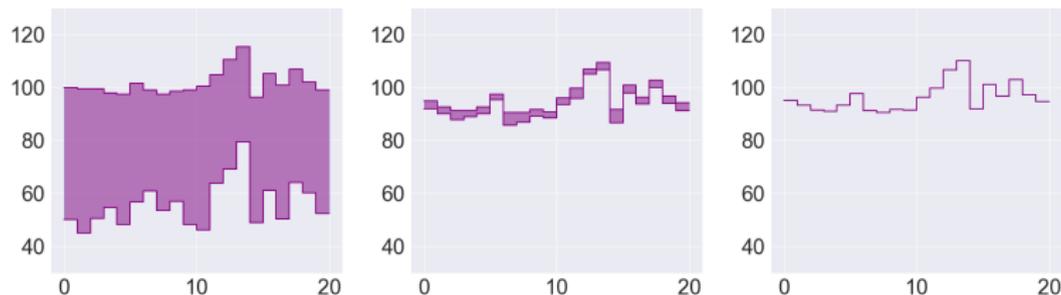


Figure: Garnet environment: evaluating π_0 , π_5 and π_{15}

Numerical results: Cartpole

- ▶ CartPole is an example of the environment with a finite action space and infinitely large state space. In this environment agent can push cart with pole on it to the left or right direction and the target is to hold the pole up as long as possible.
- ▶ Reward equals to 1 is gain every time until failing or the end of episode.
- ▶ We apply normally distributed random variable (additional randomness) to the angle.

Numerical results: Cartpole

- ▶ LD(Linear Deterministic) policy can be expressed as $I\{3 \cdot \theta + \dot{\theta} > 0\}$, where θ is an angle between pole and normal to cart.

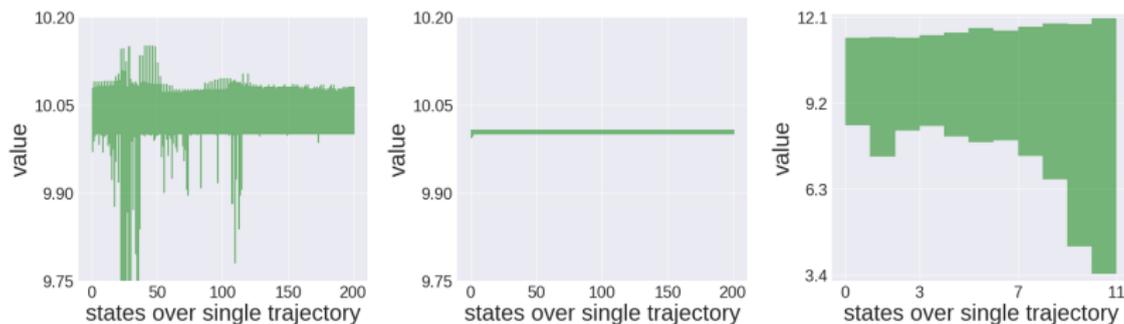


Figure: CartPole environment: evaluating A2C, deterministic policy and random policy.

UVIP: reminder

Require: $V^\pi, \widehat{V}_0^{\text{up}}, \gamma, \varepsilon, M_1, M_2$

Ensure: V^{up}

for $x \in X, a \in A$ **do**

$$\bar{V}(x, a) = M_1^{-1} \sum_{i=1}^{M_1} V^\pi(Y_i^{x,a}), \quad Y_i^{x,a} \sim P^a(\cdot|x)$$

for $y \in X$ **do**

$$\Phi_\pi^{x,a}(y) = V^\pi(y) - \bar{V}(x, a)$$

end for

end for

$k = 1$

while $\|\widehat{V}_k^{\text{up}} - \widehat{V}_{k-1}^{\text{up}}\|_X > \varepsilon$ **do**

for $x \in X$ **do**

$$\widehat{V}_{k+1}^{\text{up}}(x) = M_2^{-1} \sum_{i=1}^{M_2} [\max_a \{r^a(x) + \gamma(\widehat{V}_k^{\text{up}}(Y_i^{x,a}) - \Phi_\pi^{x,a}(Y_i^{x,a}))\}],$$
$$Y_i^{x,a} \sim P^a(\cdot|x)$$

end for

$k = k + 1$

end while

$V^{\text{up}} = \widehat{V}_k^{\text{up}}$

Further research directions: online method

- ▶ Instead of estimating $P^a V(x)$ with Monte-Carlo method, consider estimating P^a online;
- ▶ Use $\hat{P}^a V(x)$ instead of $\bar{V}(x, a)$ where

$$\hat{P}^a(x'|x) = \frac{N(x', a, x)}{N(a, x)}, x, x' \in X, a \in A;$$

- ▶ $N(a, x)$ - number of visits of pair (a, x) ,
- ▶ $N(x', a, x)$ - number of visits to x' by taking action a at state x .

UVIP: reminder

Algorithm 2: UVIP: online modification

Require: $V^\pi, \widehat{V}_0^{\text{up}}, \gamma, \varepsilon, M_1, M_2$

Ensure: V^{up}

$k = 1$

while $\|\widehat{V}_k^{\text{up}} - \widehat{V}_{k-1}^{\text{up}}\|_X > \varepsilon$ **do**

for $x \in X$ **do**

$\widehat{V}_{k+1}^{\text{up}}(x) =$

$M_2^{-1} \sum_{i=1}^{M_2} [\max_a \{r^a(x) + \gamma(\widehat{V}_k^{\text{up}}(Y_i^{x,a}) - V^\pi(Y_i^{x,a}) + \widehat{P}^a V^\pi(x))\}],$

$Y_i^{x,a} \sim P^a(\cdot|x)$

end for

$k = k + 1$

end while

$V^{\text{up}} = \widehat{V}_k^{\text{up}}$

Online method: Garnet

- Consider the sequence of policies $(\pi_k)_{k=0}^{15}$, obtained via the value iteration procedure. We evaluate policies π_k by computing $\Delta_{\pi,k}^{\text{up}}(x)$ for some time steps k .

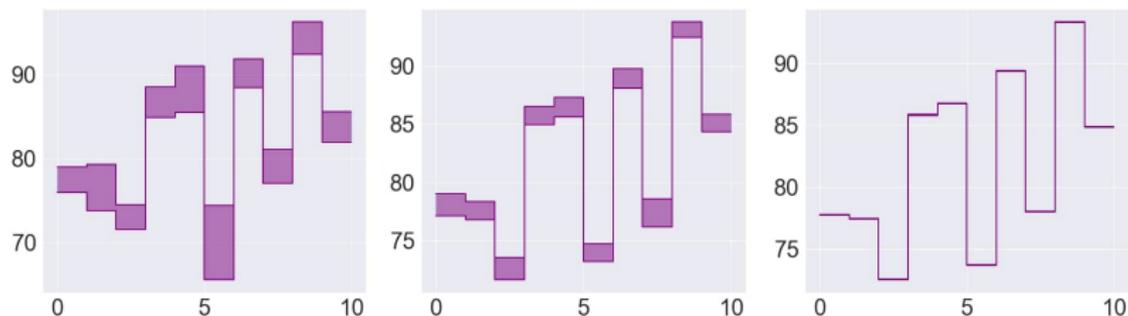


Figure: Garnet environment: evaluating π_0 , π_5 and π_{15} with an online method

Extensions: real-time UVIP

Consider now the extensions to finite-horizon case with episodes of length H .

Algorithm 3: Real-Time UVIP

Initialize: $V_t^{0,\text{up}}$.

for $k = 1, 2, \dots$ **do**

Initialize X_1^k

for $t = 0, \dots, H - 1$ **do**

$A_t^k \in \arg \max_{a \in A} \{r^a(X_t^k) + P^a V_{t+1}^{\text{up},k-1}(X_t^k)\}$

Act with A_t^k and observe X_{t+1}^k

$V_t^{\text{up},k}(X_t^k) =$

$\mathbb{E}[\max_{a \in A} \{r^a(X_t^k) + V_{t+1}^{\text{up},k-1}(Y_{t+1}^{a,k}) - V^\pi(Y_{t+1}^{a,k}) + (P^a V^\pi)(X_t^k)\}],$

$Y_{t+1}^{a,k} \sim P^a(\cdot | X_t^k).$

end for

end for

In the algorithm above $P^a V^\pi(x)$ can be replaced by $\hat{P}^a V^\pi(x)$ and the outer expectation can be replaced by its Monte-Carlo estimate.

Paper available at <https://arxiv.org/abs/2105.02135>

Thank you!

References I

- Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for reinforcement learning. 70:263–272, 06–11 Aug 2017. URL <http://proceedings.mlr.press/v70/azar17a.html>.
- Denis Belomestny and John Schoenmakers. *Advanced Simulation-Based Methods for Optimal Stopping and Control: With Applications in Finance*. Springer, 2018.
- Dimitri P Bertsekas and Steven E Shreve. Stochastic optimal control, volume 139 of mathematics in science and engineering, 1978.
- Chi Jin, Zeyuan Allen-Zhu, Sebastien Bubeck, and Michael I Jordan. Is q-learning provably efficient? 31, 2018. URL <https://proceedings.neurips.cc/paper/2018/file/d3b1fb02964aa64e257f9f26a31f72cf-Paper.pdf>.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- L. Rogers. Pathwise stochastic optimal control. *SIAM J. Control and Optimization*, 46:1116–1132, 01 2007. doi: 10.1137/050642885.
- Ibrahim El Shar and Daniel Jiang. Lookahead-bounded q-learning. 119:8665–8675, 13–18 Jul 2020. URL <http://proceedings.mlr.press/v119/shar20a.html>.
- Richard Sutton. Learning to predict by the method of temporal differences. *Machine Learning*, 3: 9–44, 08 1988. doi: 10.1007/BF00115009.
- J. N. Tsitsiklis and B. Van Roy. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5):674–690, May 1997. ISSN 2334-3303. doi: 10.1109/9.580874.