

Long-run Average Reward for Markov Decision Processes

Based on a paper at CAV 2017

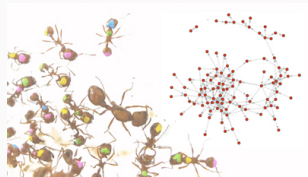
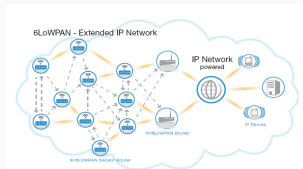
Pranav Ashok¹, Krishnendu Chatterjee², Przemysław Daca²,
Jan Křetínský¹ and Tobias Meggendorfer¹

August 9, 2017

¹Technical University of Munich, Germany

²IST Austria

Motivation



Markov Decision Processes (MDPs) are a standard model for describing systems which display probabilistic as well as non-deterministic behaviour.

Open Problem

Long-run Average Reward or Mean-payoff using Value Iteration (VI)

- VI approaches exist for subclasses, but not for general MDPs
- Other existing approaches: Linear Programming (LP) and Strategy Iteration (SI)¹

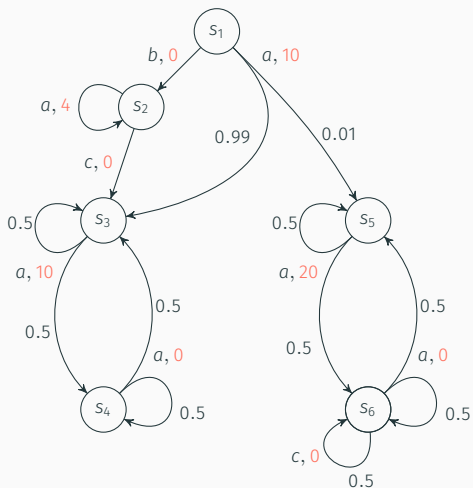
¹Martin L. Puterman, Markov Decision Processes, 1994.

Contributions

- **Disprove** conjectured stopping criterion for VI²
- **General solution** using VI
- Improve performance using ideas from **Machine Learning**

²Not covered in this talk

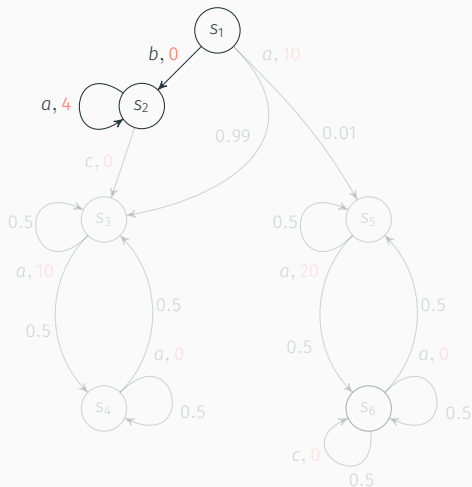
Markov Decision Process (MDPs)



Strategy

A strategy (or policy) gives the action to be taken at every state.

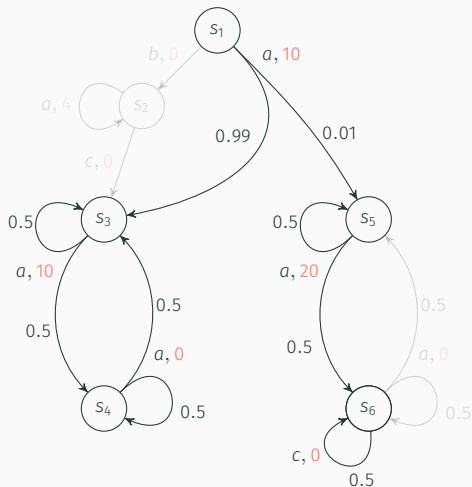
$$\pi := \{s_1 \mapsto b, s_2 \mapsto a\}$$



Strategy

A strategy (or policy) gives the action to be taken at every state.

$$\pi := \{s_1 \mapsto a, \dots, s_6 \mapsto a\}$$



Mean-payoff or Long-run Average Reward

$$\rho = 15 \ 35 \ 20 \ 50 \ 0 \ 10 \ 0 \ 10 \ 0 \ 10 \ 0 \dots$$

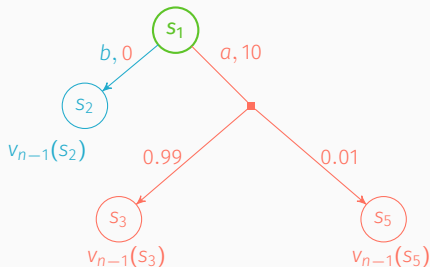
Then, n-step average reward is given by

$$MP_n(\rho) := \frac{1}{n} \cdot \sum_{i=1}^n \rho_i$$

$$MP := \sup_{\pi} \liminf_{n \rightarrow \infty} \mathbb{E}^{\pi} [MP_n(\rho)]$$

VI for Total Rewards (Bellman Equation)

$$v_n(s_1) = \max \{ 0 + v_{n-1}(s_2), 10 + (0.99 \cdot v_{n-1}(s_3) + 0.01 \cdot v_{n-1}(s_5)) \}$$

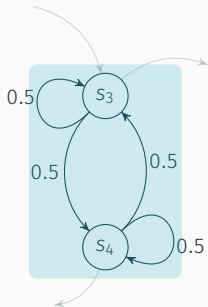


Average reward³

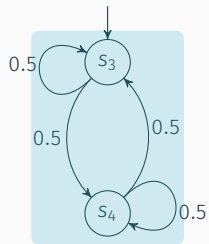
$$\lim_{n \rightarrow \infty} \frac{v_n(s)}{n} \approx v_n(s) - v_{n-1}(s)$$

³After making the MDP aperiodic

Towards a General VI: Maximal End Components (MECs)



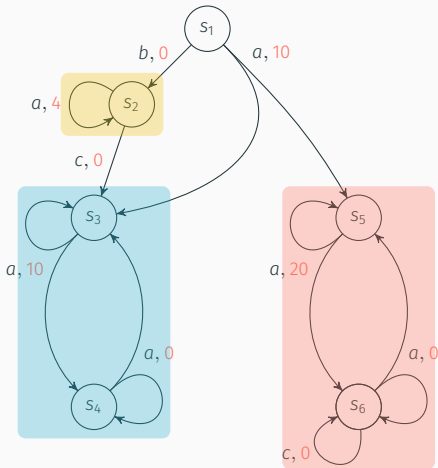
An example MEC



Communicating MDP

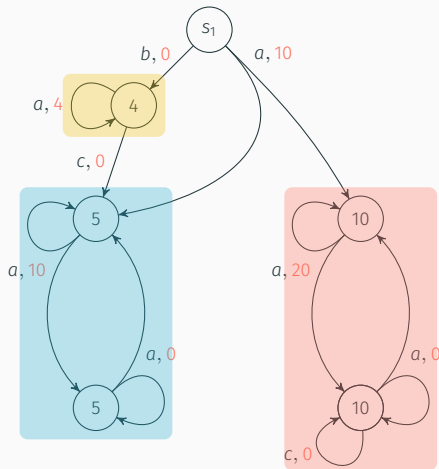
VI ✓

Towards a General VI: Step 1 – MEC-Decomposition



Find all MECs

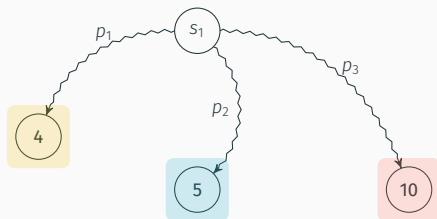
Towards a General VI: Step 1 – MEC-Decomposition



Find all MECs and **run VI on them until ϵ -convergence**

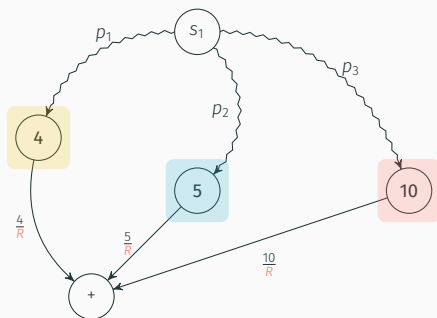
Towards a General VI: Step 2 – Weighted Reachability

$$\text{Max. mean-payoff} = \sup_{\pi} p_1 \cdot 4 + p_2 \cdot 5 + p_3 \cdot 10$$



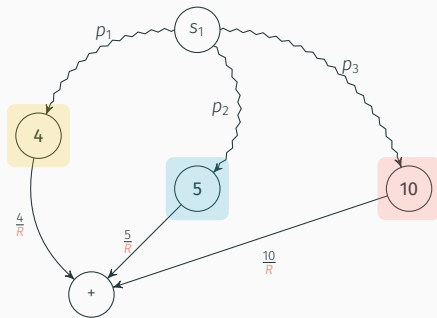
Towards a General VI: Step 2 – Weighted Reachability

$$\frac{\text{Max. mean-payoff}}{R} = \sup_{\pi} p_1 \frac{4}{R} + p_2 \frac{5}{R} + p_3 \frac{10}{R}$$



Towards a General VI: Step 2 – Weighted Reachability

$$\frac{\text{Max. mean-payoff}}{R} = \sup_{\pi} p_1 \frac{4}{R} + p_2 \frac{5}{R} + p_3 \frac{10}{R}$$



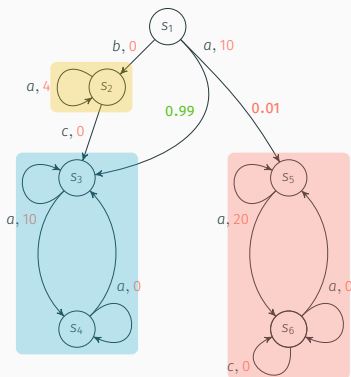
Mean-payoff reduced to reachability: $P_{max}(\diamond +)$

Limitations of this method

- Whole state space is being explored
- Finding MECs on large state spaces is computationally expensive

Improvement: avoid full state-space exploration

Idea: let sampling guide us to the “important” regions



Contribution of red region to MP is potentially low

Improvement: guarantees through sampling

- Existing: BRTDP approach for reachability^{4,5}
- Repeatedly samples paths from initial state
- Back-propagates values along the path using VI operator

⁴Bounded Real-Time Dynamic Programming, McMahan et. al., ICML '05

⁵Verification of Markov Decision Processes using Learning Algorithms, Brazdil et. al., ATVA '14

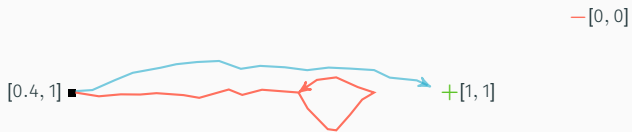
Improvement: BRTDP in action



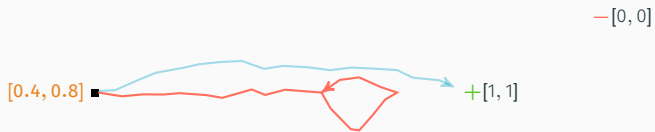
Improvement: BRTDP in action



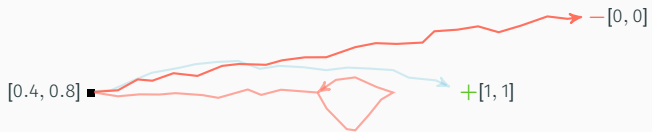
Improvement: BRTDP in action



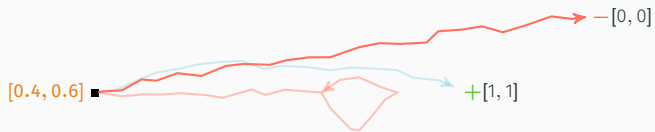
Improvement: BRTDP in action



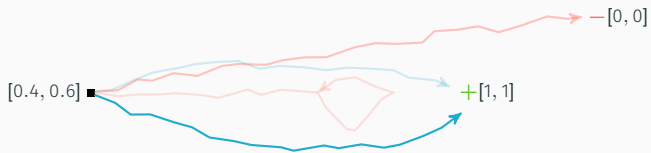
Improvement: BRTDP in action



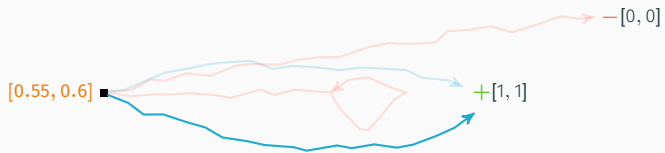
Improvement: BRTDP in action



Improvement: BRTDP in action



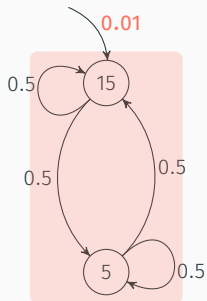
Improvement: BRTDP in action



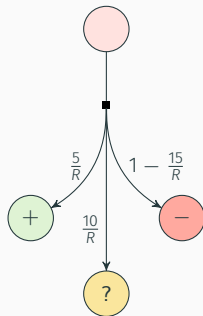
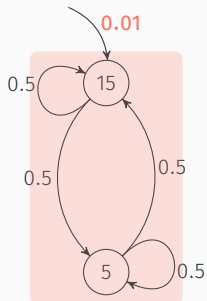
Improvement: Algorithm

1. Run BRTDP in search of the $+$ state
2. Repeatedly collapse MECs amongst the states explored so far
3. While collapsing, compute their values and add special transition to $+$ and $-$ states

Improvement: can do even more...



Improvement: can do even more...



Collapse MEC into a single state and add a special action

Final Algorithm: On-demand Value Iteration (ODV)

1. Sample paths like in BRTDP
2. When MECs are detected, collapse them but...
3. Don't compute MEC value until ϵ -convergence
4. Add transition with probability $\propto (U-L)$ to ? state
5. Refine value of MEC only when ? encountered

We saw two methods which can be used to obtain mean-payoff with guarantees

1. Collapse MECs, add transitions to $+/-$ states, run reachability
2. **ODV**: Run sampling, collapse on-the-go, refine MEC values on-demand

Benchmarks

Model	States	MECs	LP ¹	MEC-VI ²
virus	809	1	0.19	0.05
cs_nfail4	960	176	0.7	0.18
investor	6 688	837	2.8	0.51
phil-nofair5	93 068	1	TO	6.67
rabin4	668 836	1	TO	112.38

1. MultiGain, Brazdil et. al. 2015.
2. MEC-VI: Straightforward conversion to reachability, then VI

Benchmarks

On-demand VI better by orders of magnitude depending on topology

Model	States	MEC-VI	ODV	ODV States	ODV MECs
zeroconf(40,10)	3 001 911	MO	5.05	481	3
<i>avoid</i>				582	3
zeroconf(300,15)	4 730 203	MO	16.6	873	3
<i>avoid</i>				5 434	3
sensors(2)	7 860	18.9	20.1	3 281	917
sensors(3)	77 766	2293.0	37.0	10 941	2 301

