Based on [BCCF+14]:
Verification of Markov decision processes using learning algorithms, ATVA 2014.

# I/ MDPs



**MDP** $\mathcal{M} = (S, A, \delta)$

- $S$: finite set of states
- $A$: " " " actions
- $\delta$: $S \times A \to \mathcal{D}(S)$ partial probabilistic transit⁰ funct⁰

**Strategy** $\sigma \in \Sigma$ : $(S \times A)^* S \to \mathcal{D}(A)$ s.t.

$$\sigma(\wedge, \alpha) \text{ is defined}$$
for each $\alpha \in \text{Supp}(\sigma(\rho))$
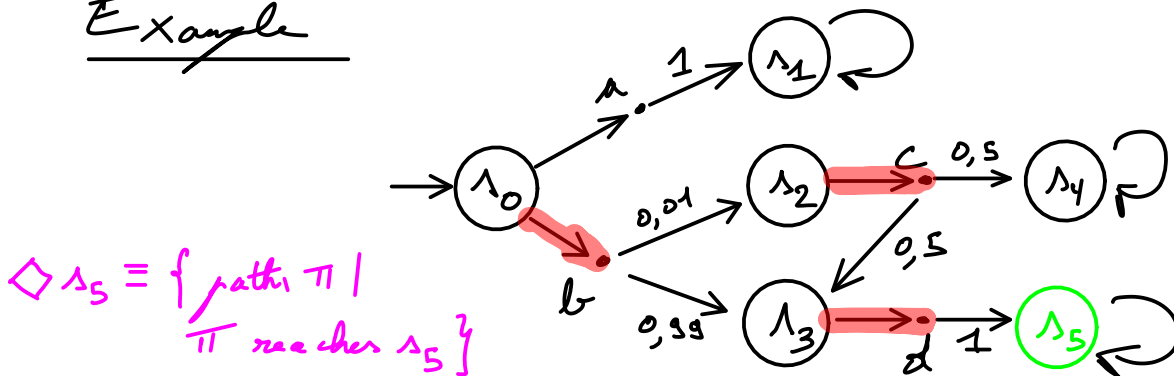
<u>Rmq</u> Strategies may use memory and randomness.

Given an MDP $\mathcal{M}$, a state $\wedge$, a strategy $\sigma$, we obtain a Markov chain (MC), hence a probability measure over paths: $\mathbb{P}^\sigma_{\mathcal{M}, \wedge}$.

# II | Reachability

The goal is to reach $T \subseteq S$ and we want a strategy $\sigma^*$ that maximizes the probability to do so.

**Rmk** Lots of problems on MDPs can be reduced to reachability problems!

## Example



$\Diamond s_5 \equiv \{ \text{paths } \pi \mid \pi \text{ reaches } s_5 \}$

$$\sup_{\sigma} \mathbb{P}^{\sigma}_{s_0} [\Diamond s_5] = 1 - 0,01 \cdot 0,5 = 0,995$$

> **Thm**
> Pure memoryless strategies suffice for reachability.

$\hookrightarrow$ We may replace $\sup_{\sigma \in \Sigma}$ by $\max_{\sigma \in \Sigma^{PM}}$ !

$\implies$ Cornerstone of the following algorithms.

# III | Classical algorithms

## a) Linear Program

Vector $(x_s)_{s \in S}$ with $x_s = \max_{\sigma} \mathbb{P}^{\sigma}_s[\lozenge T]$ is the unique solut° of the LP:

- If $s \in T$, then $x_s = 1$.
- If $s \not\models \exists \lozenge T$, then $x_s = 0$.

  $\underbrace{\qquad\qquad}$
  $s$ not connected to $T$

- Else, $0 \leq x_s \leq 1$
  and for all $\alpha \in A(s)$ (act° enabled in $s$)
  $$x_s \geq \sum_{t \in S} \delta(s, \alpha, t) \cdot x_t \quad (\ast)$$

where $\sum_{s \in S} x_s$ is ==minimal==. $(\ast\ast)$

Intuit°
$(\ast)$ Otherwise we could increase $x_s$ by changing $\alpha$
$(\ast\ast)$ We need the smaller fixed point
  ( $x_s = 1$ is always a solut° but not the one
  we want)

$\implies$ Problem $\in P$. But in practice, LP does not scale well.

# b) Value & strategy iteration

VI : approximation technique for values $x_s$

1. Backward reachability to determine
$$\{ s \mid s \models \exists \Diamond T \} = \{ s \mid x_s > 0 \} = Pre^*(T)$$

2. For $s \in Pre^*(T) \setminus T$ :
$$x_s = \lim_{m \to \infty} x_s^{(m)}$$

where $x_s^{(0)} = 0$
and $x_s^{(m+1)} = \max \left\{ \sum_{t \in S} \delta(s, \alpha, t) \cdot x_t^{(m)} \mid \alpha \in A(s) \right\}$

*When $=$, Bellman equation*

Stopping criterion:
(i) Naïve : stop when $\max_{s \in S} | x_s^{(m+1)} - x_s^{(m)} | < \varepsilon$ for some $\varepsilon$.
$\hookrightarrow$ Fails in some cases

(ii) Also compute an upper bound and evaluate the difference.
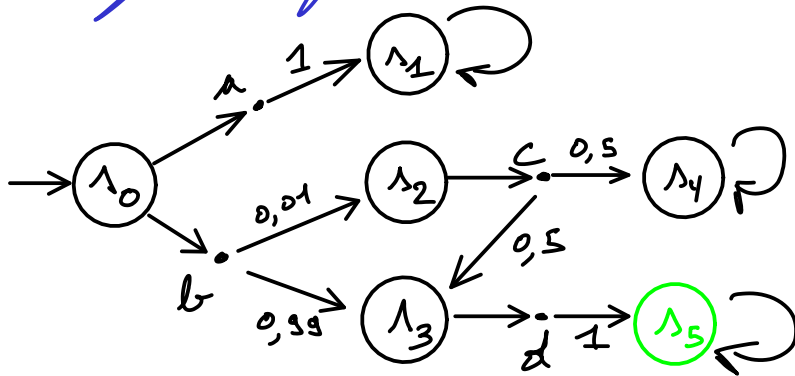
**Thm** VI converges monotonically.

$\longrightarrow$ Exponential in the worst-case but the most efficient technique in practice.

**Rmk** VI does not give an $(\varepsilon-)$ optimal scheduler directly and obtaining it from the values may be costly.

# Example of VI



| Iteration | $\Lambda_0$ | $\Lambda_1$ | $\Lambda_2$ | $\Lambda_3$ | $\Lambda_4$ | $\Lambda_5$ |
|-----------|------|------|------|------|------|------|
| 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 1 | 0 | 0 | 0 | 1 | 0 | 1 |
| 2 | 0,99 | 0 | 0,5 | 1 | 0 | 1 |
| 3 | 0,995 | 0 | 0,5 | 1 | 0 | 1 |

$\longrightarrow$ Here we reach the limit in a finite time because no loop $\longrightarrow$ not true in general.

**SI**: we want to compute actual strategies, not merely values

We start as before with $x_s^{(0)}$ for $s \in \text{Pre}^*(T) \setminus T$. Then we loop:

(i) $\sigma^{(m+1)}(s) = \underset{\alpha \in A(s)}{\arg\max} \left\{ \sum_{t \in S} \delta(s, \alpha, t) \cdot x_t^{(m)} \right\}$

(ii) Evaluate $x_s^{(m+1)} = \mathbb{P}_s^{\sigma^{(m+1)}}[\Diamond T]$ either using VI (approx.) or linear equation system (exact but slower).

Stopping criterion: when $\sigma^{(m+1)} = \sigma^{(m)}$

**Thm**
SI converges monotonically.

$\longrightarrow$ Also exponential in the worst-case. Slower than VI but more precise.

**Comparison**

Similar log for both VI and SI

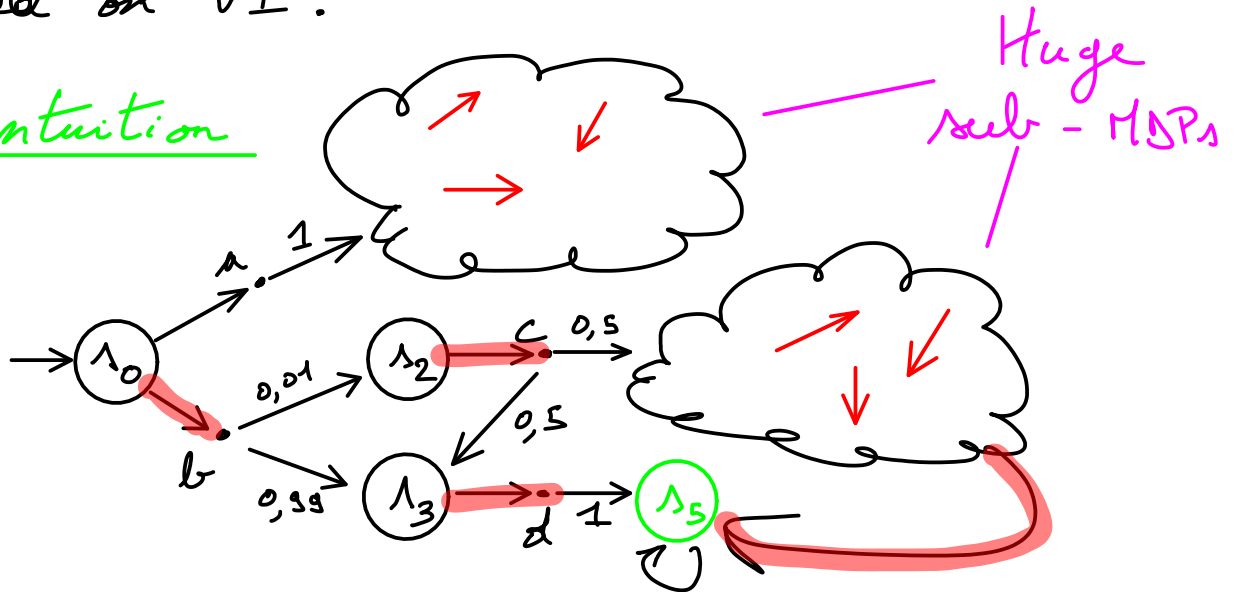1) Choose best actions based on current approx.

2) Update current approx.

$\hookrightarrow$ In VI, update over one step using previous approx

$\hookrightarrow$ In SI, "exact" computation based on the chosen actions.

Using learning: BRTDP approach

Many different variants: I consider one
based on VI.

**Intuition**



We do not care <u>at all</u> about the
first cloud, and we can estimate the
optimal probability up to $\varepsilon = 0,005$
without looking at the second cloud.

$\Longrightarrow$ We want to identify which
parts of the MDP are really important
and focus on them.

$\Longrightarrow$ Will speed up the search for
good strategies but also permit
to highlight meaningful actions, hence
granting simple strategies

E.g., here $\sigma(s) = \begin{cases} b \text{ if possible} \\ \text{random otherwise} \end{cases}$

already yields $\mathbb{P}^{\sigma}[\lozenge T] \geq 0,99$

$\Longrightarrow$ Related research: strategies as decision trees.

: strongly - connected sub - MDP

(i.e., no probability leak).

$\rightsquigarrow$ similar to BSCCs for MCs.

==Simpler case :== - trivial ECs 🙂 & ☹
                                    target      sink

- initial state $\overline{\lambda}$

We want $\max\limits_{\sigma \in \Sigma^{PM}} \mathbb{P}^{\sigma}_{\mathcal{M}, \overline{\lambda}} [\Diamond 🙂] = x_{\overline{\lambda}}$

$\underline{Rmq}$ : we will approximate $x_{\overline{\lambda}, \alpha}$

$$= \sum_{t \in S} \delta(\overline{\lambda}, \alpha, t) \cdot x_t$$

but $x_{\overline{\lambda}} = \max\limits_{\alpha \in A(\lambda)} x_{\overline{\lambda}, \alpha}$

==Intuition==: we will approximate $x_{\overline{\lambda}, \alpha}$
using upper and lower bounds
$\hookrightarrow U(\overline{\lambda}, \alpha) \longrightarrow L(\overline{\lambda}, \alpha)$

$$U(\cdot,\cdot) \leftarrow 1, \quad L(\cdot,\cdot) \leftarrow 0$$
$$L(\smiley,\cdot) \leftarrow 1, \quad U(\frownie,\cdot) \leftarrow 0$$

$\rbrace$ *Trivial bounds*

REPEAT
   $\rho \leftarrow \overline{\lambda}$
   REPEAT
      $\alpha \leftarrow$ sampled uniformaly from
        $\underset{\alpha \in A(last(\rho))}{arg\ max} U(last(\rho), \alpha)$

*last state of prefix $\rho$*

      $\hookrightarrow$ *random action maximizing $U$*
        $\Rightarrow$ *Focus on "good" strategies*

      $\lambda \leftarrow$ sampled according to $\delta(last(\rho), \alpha)$
      $\rho \leftarrow \rho \cdot \alpha \cdot \lambda$
   UNTIL $\lambda \in \{\smiley, \frownie\}$   $\longrightarrow$ *Path $\rho$ ends in one of the ECs*

   REPEAT
      $\lambda' \leftarrow pop(\rho)$
      $\alpha \leftarrow pop(\rho)$
      $\lambda \leftarrow last(\rho)$
      UPDATE$(\lambda, \alpha)$
   UNTIL $\rho = \overline{\lambda}$

$\rbrace$ *For each visited transit°, update $U$ and $L$ (backwards)*

UNTIL $U(\overline{\lambda}) - L(\overline{\lambda}) < \varepsilon$

$\hookrightarrow$ *Approx. is close enough.*

<u>with</u>
$$U(\lambda) = \underset{\alpha \in A(\lambda)}{max}\ U(\lambda, \alpha)$$
$$L(\lambda) = \underset{\alpha \in A(\lambda)}{max}\ L(\lambda, \alpha)$$

*Explore phase*

*Update phase*

UPDATE $(s, \alpha)$

$$U(s, \alpha) = \sum_{s' \in S} \delta(s, \alpha, s') \cdot U(s')$$

$$L(s, \alpha) = \sum_{s' \in S} \delta(s, \alpha, s') \cdot L(s')$$

**Rmq 1** This works if the MDP is known.
If $\delta$ is not known but can be sampled,
UPDATE needs to be adapted
$\Longrightarrow$ It gives a PAC algorithm.
probably approximately correct

**Rmq 2** The general case (non-trivial ECs)
is more complex.
Intuition:
— identify ECs from long enough simulations
— contract them on the fly

# V | Conclusion

1) Reachability in MDPs is a key problem.
2) Textbook solut°: LP
$\longrightarrow$ not efficient in practice
3) Classical technique: value / strategy iterat°
$\longrightarrow$ works well in practice
4) Mix with learning
$\longrightarrow$ several variants (heuristics)
$\longrightarrow$ tremendous gains: state space $\times 10^{-3}$,
time $\times 10^{-2}$
5) Quantitative extensions
$\longrightarrow$ e.g., mean-payoff