

Modeling and Verification of Probabilistic Systems

Joost-Pieter Katoen

Lehrstuhl für Informatik 2
Software Modeling and Verification Group

<http://moves.rwth-aachen.de/teaching/ss-14/movep14/>

May 22, 2014

Markov decision process (MDP)

Markov decision processes

- ▶ In MDPs, **both** nondeterministic and probabilistic choices coexist.
- ▶ MDPs are transition systems in which in any state a nondeterministic choice between probability distributions exists.
- ▶ Once a probability distribution has been chosen nondeterministically, the next state is selected probabilistically—as in DTMCs.
- ▶ Any MC is thus an MDP in which in any state the probability distribution is uniquely determined.

Randomized distributed algorithms are typically appropriately modeled by MDPs, as probabilities affect just a small part of the algorithm and nondeterminism is used to model concurrency between processes by means of interleaving.

Overview

- 1 Markov Decision Processes
- 2 Probabilities in MDPs
- 3 Policies
 - Positional policies
 - Finite-memory policies
- 4 Reachability probabilities
 - Mathematical characterisation
 - Value iteration
 - Linear programming
 - Policy iteration
- 5 Summary

Markov decision process (MDP)

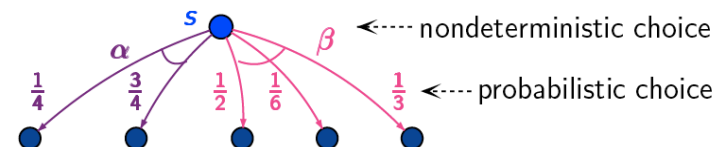
Markov decision process

An MDP \mathcal{M} is a tuple $(S, Act, \mathbf{P}, \ell_{init}, AP, L)$ where

- ▶ S is a countable set of states with initial distribution $\ell_{init} : S \rightarrow [0, 1]$
- ▶ Act is a finite set of actions
- ▶ $\mathbf{P} : S \times Act \times S \rightarrow [0, 1]$, transition probability function such that:

$$\text{for all } s \in S \text{ and } \alpha \in Act : \sum_{s' \in S} \mathbf{P}(s, \alpha, s') \in \{0, 1\}$$

- ▶ AP is a set of atomic propositions and labeling $L : S \rightarrow 2^{AP}$.



Markov decision process (MDP)

Markov decision process

An MDP \mathcal{M} is a tuple $(S, Act, \mathbf{P}, \iota_{\text{init}}, AP, L)$ where

- ▶ $S, \iota_{\text{init}} : S \rightarrow [0, 1]$, AP and L are as before, i.e., as for DTMCs, and
- ▶ Act is a finite set of actions
- ▶ $\mathbf{P} : S \times Act \times S \rightarrow [0, 1]$, transition probability function such that:

$$\text{for all } s \in S \text{ and } \alpha \in Act : \sum_{s' \in S} \mathbf{P}(s, \alpha, s') \in \{0, 1\}$$

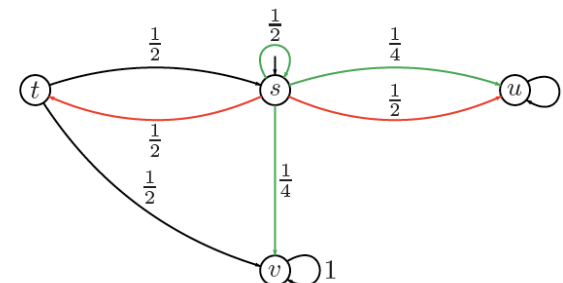
Enabled actions

Let $Act(s) = \{\alpha \in Act \mid \exists s' \in S. \mathbf{P}(s, \alpha, s') > 0\}$ be the set of enabled actions in state s . We require $Act(s) \neq \emptyset$ for any state s .

Overview

- 1 Markov Decision Processes
- 2 Probabilities in MDPs
- 3 Policies
 - Positional policies
 - Finite-memory policies
- 4 Reachability probabilities
 - Mathematical characterisation
 - Value iteration
 - Linear programming
 - Policy iteration
- 5 Summary

An example MDP



- ▶ Initial distribution: $\iota_{\text{init}}(s) = 1$ and $\iota_{\text{init}}(t) = \iota_{\text{init}}(u) = \iota_{\text{init}}(v) = 0$
- ▶ Set of enabled actions in state s is $Act(s) = \{\alpha, \beta\}$ where
 - ▶ $\mathbf{P}(s, \alpha, s) = \frac{1}{2}$, $\mathbf{P}(s, \alpha, t) = 0$ and $\mathbf{P}(s, \alpha, u) = \mathbf{P}(s, \alpha, v) = \frac{1}{4}$
 - ▶ $\mathbf{P}(s, \beta, s) = \mathbf{P}(s, \beta, v) = 0$, and $\mathbf{P}(s, \beta, t) = \mathbf{P}(s, \beta, u) = \frac{1}{2}$
- ▶ $Act(t) = \{\alpha\}$ with $\mathbf{P}(t, \alpha, s) = \mathbf{P}(t, \alpha, u) = \frac{1}{2}$ and 0 otherwise

Paths in an MDP

State graph

The *state graph* of MDP \mathcal{M} is a digraph $G = (V, E)$ with V are the states of M , and $(s, s') \in E$ iff $\mathbf{P}(s, \alpha, s') > 0$ for some $\alpha \in Act$.

Paths

An infinite *path* in an MDP $\mathcal{M} = (S, Act, \mathbf{P}, \iota_{\text{init}}, AP, L)$ is an infinite sequence $s_0 \alpha_1 s_1 \alpha_2 s_2 \alpha_3 \dots \in (S \times Act)^\omega$, written as

$$\pi = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots,$$

such that $\mathbf{P}(s_i, \alpha_{i+1}, s_{i+1}) > 0$ for all $i \geq 0$. Any finite prefix of π that ends in a state is a *finite path*.

Let $Paths(\mathcal{M})$ denote the set of paths in \mathcal{M} , and $Paths^*(\mathcal{M})$ the set of finite prefixes thereof.

Overview

- 1 Markov Decision Processes
- 2 Probabilities in MDPs
- 3 Policies
 - Positional policies
 - Finite-memory policies
- 4 Reachability probabilities
 - Mathematical characterisation
 - Value iteration
 - Linear programming
 - Policy iteration
- 5 Summary

Induced DTMC of an MDP by a policy

DTMC of an MDP induced by a policy

Let $\mathcal{M} = (S, Act, \mathbf{P}, \iota_{init}, AP, L)$ be an MDP and \mathfrak{G} a policy on \mathcal{M} . The DTMC *induced by* \mathfrak{G} , denoted $\mathcal{M}_{\mathfrak{G}}$, is given by

$$\mathcal{M}_{\mathfrak{G}} = (S^+, \mathbf{P}_{\mathfrak{G}}, \iota_{init}, AP, L')$$

where for $\sigma = s_0 s_1 \dots s_n$: $\mathbf{P}_{\mathfrak{G}}(\sigma, \sigma s_{n+1}) = \mathbf{P}(s_n, \mathfrak{G}(\sigma), s_{n+1})$ and $L'(\sigma) = L(s_n)$.

$\mathcal{M}_{\mathfrak{G}}$ is infinite, even if the MDP \mathcal{M} is finite. Since policy \mathfrak{G} might select different actions for finite paths that end in the same state s , a policy as defined above is also referred to as *history-dependent*.

Policies

Policy

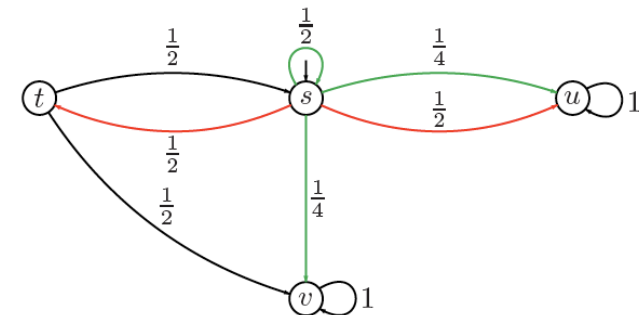
Let $\mathcal{M} = (S, Act, \mathbf{P}, \iota_{init}, AP, L)$ be an MDP. A *policy* for \mathcal{M} is a function $\mathfrak{G} : S^+ \rightarrow Act$ such that $\mathfrak{G}(s_0 s_1 \dots s_n) \in Act(s_n)$ for all $s_0 s_1 \dots s_n \in S^+$.

The path

$$\pi = s_0 \xrightarrow{\alpha_1} s_1 \xrightarrow{\alpha_2} s_2 \xrightarrow{\alpha_3} \dots$$

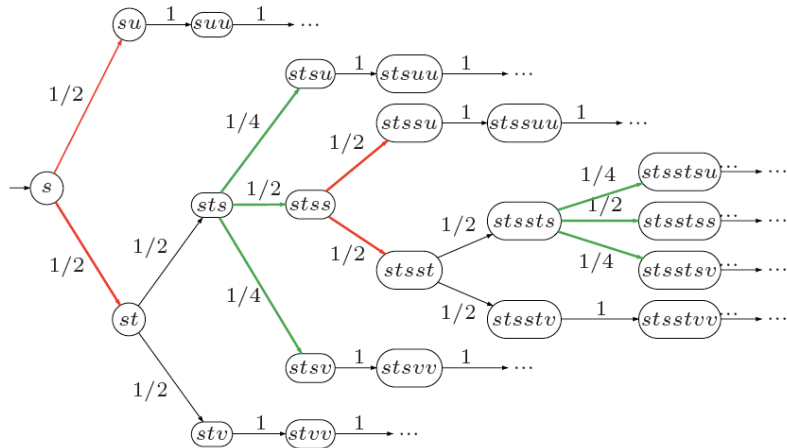
is called a \mathfrak{G} -path if $\alpha_i = \mathfrak{G}(s_0 \dots s_{i-1})$ for all $i > 0$.

Example MDP



Consider a policy that alternates between selecting **red** and **green**, starting with **red**.

Example induced DTMC



Induced DTMC for a policy that alternates between selecting **red** and **green**.

Positional policy

Positional policy

Let \mathcal{M} be an MDP with state space S . Policy \mathfrak{G} on \mathcal{M} is *positional* (or: *memoryless*) iff for each sequence $s_0 s_1 \dots s_n$ and $t_0 t_1 \dots t_m \in S^+$ with $s_n = t_m$:

$$\mathfrak{G}(s_0 s_1 \dots s_n) = \mathfrak{G}(t_0 t_1 \dots t_m).$$

In this case, \mathfrak{G} can be viewed as a function $\mathfrak{G} : S \rightarrow Act$.

Policy \mathfrak{G} is positional if it always selects the same action in a given state. This choice is independent of what has happened in the history, i.e., which path led to the current state.

Probability measure on MDP

Probability measure on MDP

Let $Pr_{\mathfrak{G}}^{\mathcal{M}}$, or simply $Pr^{\mathfrak{G}}$, denote the probability measure $Pr^{\mathcal{M}_{\mathfrak{G}}}$ associated with the DTMC $\mathcal{M}_{\mathfrak{G}}$.

This measure is the basis for associating probabilities with events in the MDP \mathcal{M} . Let, e.g., $P \subseteq (2^{AP})^{\omega}$ be an ω -regular property. Then $Pr^{\mathfrak{G}}(P)$ is defined as:

$$Pr^{\mathfrak{G}}(P) = Pr^{\mathcal{M}_{\mathfrak{G}}}(P) = Pr_{\mathcal{M}_{\mathfrak{G}}}\{\pi \in Paths(\mathcal{M}_{\mathfrak{G}}) \mid trace(\pi) \in P\}.$$

Similarly, for fixed state s of \mathcal{M} , which is considered as the unique starting state,

$$Pr^{\mathfrak{G}}(s \models P) = Pr_s^{\mathcal{M}_{\mathfrak{G}}}\{\pi \in Paths(s) \mid trace(\pi) \in P\}$$

where we identify the paths in $\mathcal{M}_{\mathfrak{G}}$ with the corresponding \mathfrak{G} -paths in \mathcal{M} .

Finite-memory policies

- ▶ *Finite-memory policies* (shortly: fm-policies) are a generalisation of positional policies.
- ▶ The behavior of an fm-policy is described by a deterministic finite automaton (DFA).
- ▶ The selection of the action to be performed in the MDP \mathcal{M} depends on the current state of \mathcal{M} (as before) and the current state (called *mode*) of the policy, i.e., the DFA.

Finite-memory policy

Finite-memory policy

Let \mathcal{M} be an MDP with state space S and action set Act .

A *finite-memory policy* \mathfrak{G} for \mathcal{M} is a tuple $\mathfrak{G} = (Q, act, \Delta, start)$ with:

- ▶ Q is a finite set of **modes**,
- ▶ $\Delta : Q \times S \rightarrow Q$ is the **transition function**,
- ▶ $act : Q \times S \rightarrow Act$ is a function that selects an action $act(q, s) \in Act(s)$ for any mode $q \in Q$ and state $s \in S$ of \mathcal{M} ,
- ▶ $start : S \rightarrow Q$ is a function that selects a **starting mode** for state $s \in S$.

Finite-memory policies

Relation fm-policy to definition policy

An fm-policy $\mathfrak{G} = (Q, act, \Delta, start)$ is identified with policy, $\mathfrak{G}' : Paths^* \rightarrow Act$ which is defined as follows.

1. For the starting state s_0 , let $\mathfrak{G}'(s_0) = act(start(s_0), s_0)$.
2. For path fragment $\hat{\pi} = s_0 s_1 \dots s_n$ let

$$\mathfrak{G}'(\hat{\pi}) = act(q_n, s_n)$$

where $q_0 = start(s_0)$ and $q_{i+1} = \Delta(q_i, s_i)$ for $0 \leq i \leq n$.

Positional policies can be considered as fm-policies with just a single mode.

An MDP under a finite-memory policy

The behavior of an MDP \mathcal{M} under fm-policy $\mathfrak{G} = (Q, act, \Delta, start)$ is:

- ▶ Initially, a starting state s_0 is randomly determined according to the initial distribution ι_{init} , i.e., $\iota_{init}(s_0) > 0$.
- ▶ The fm-policy \mathfrak{G} initializes its DFA to the mode $q_0 = start(s_0) \in Q$.
- ▶ If \mathcal{M} is in state s and the current mode of \mathfrak{G} is q , then the decision of \mathfrak{G} , i.e., the selected action, is $\alpha = act(q, s) \in Act(s)$.
- ▶ The policy changes to mode $\Delta(q, s)$, while \mathcal{M} performs the selected action α and randomly moves to the next state according to the distribution $\mathbf{P}(s, \alpha, \cdot)$.

The DTMC under an fm-policy

Remark

For fm-policy \mathfrak{G} , the DTMC $\mathcal{M}_{\mathfrak{G}}$ can be identified with a DTMC $\mathcal{M}'_{\mathfrak{G}}$, say, where the states are just pairs $\langle s, q \rangle$ where s is a state in the MDP \mathcal{M} and q a mode of \mathfrak{G} .

$\mathcal{M}'_{\mathfrak{G}}$ is the DTMC with state space $S \times Q$, labeling $L'(\langle s, q \rangle) = L(s)$, the starting distribution ι_{init} , and the transition probabilities:

$$\mathbf{P}'_{\mathfrak{G}}(\langle s, q \rangle, \langle t, p \rangle) = \mathbf{P}(s, act(q, s), t).$$

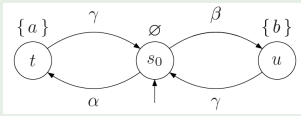
For any MDP \mathcal{M} and fm-policy \mathfrak{G} : $\mathcal{M}_{\mathfrak{G}} \sim_p \mathcal{M}'_{\mathfrak{G}}$.

Hence, if \mathcal{M} is a finite MDP, then $\mathcal{M}_{\mathfrak{G}}$ is bisimilar to the finite DTMC $\mathcal{M}'_{\mathfrak{G}}$.

Positional versus fm-policies

Positional policies are insufficient for ω -regular properties

Consider the MDP:



Positional policy \mathfrak{S}_α always chooses α in state s_0

Positional policy \mathfrak{S}_β always chooses β in state s_0 . Then:

$$Pr_{\mathfrak{S}_\alpha}(s_0 \models \diamond a \wedge \diamond b) = Pr_{\mathfrak{S}_\beta}(s_0 \models \diamond a \wedge \diamond b) = 0.$$

Now consider fm-policy $\mathfrak{S}_{\alpha\beta}$ which alternates between selecting α and β .

Then: $Pr_{\mathfrak{S}_{\alpha\beta}}(s_0 \models \diamond a \wedge \diamond b) = 1$.

Thus, the class of positional policies is insufficiently powerful to characterise minimal (or maximal) probabilities for ω -regular properties.

Other kinds of policies

- ▶ **Counting** policies that base their decision on the number of visits to a state, or the length of the history (i.e., number of visits to all states)
- ▶ **Partial-observation** policies that base their decision on the trace $L(s_0) \dots L(s_n)$ of the history $s_0 \dots s_n$.
- ▶ **Randomised** policies. This is applicable to all (deterministic) policies. For instance, a randomised positional policy $\mathfrak{S} : S \rightarrow \text{Dist}(\text{Act})$, where $\text{Dist}(X)$ is the set of probability distributions on X , such that $\mathfrak{S}(s)(\alpha) > 0$ iff $\alpha \in \text{Act}(s)$. Similar can be done for fm-policies and history-dependent policies etc..
- ▶ There is a **strict hierarchy** of policies, showing their expressiveness (black board).

Overview

- 1 Markov Decision Processes
- 2 Probabilities in MDPs
- 3 Policies
 - Positional policies
 - Finite-memory policies
- 4 **Reachability probabilities**
 - Mathematical characterisation
 - Value iteration
 - Linear programming
 - Policy iteration
- 5 Summary

Reachability probabilities

Reachability probabilities

Let \mathcal{M} be an MDP with state space S and \mathfrak{S} be a policy on \mathcal{M} . The **reachability probability** of $G \subseteq S$ from state $s \in S$ under policy \mathfrak{S} is:

$$Pr^{\mathfrak{S}}(s \models \diamond G) = Pr_s^{\mathcal{M}^{\mathfrak{S}}} \{ \pi \in \text{Paths}(s) \mid \pi \models \diamond G \}$$

Maximal and minimal reachability probabilities

The **minimal** reachability probability of $G \subseteq S$ from $s \in S$ is:

$$Pr^{\min}(s \models \diamond G) = \inf_{\mathfrak{S}} Pr^{\mathfrak{S}}(s \models \diamond G)$$

In a similar way, the **maximal** reachability probability of $G \subseteq S$ is:

$$Pr^{\max}(s \models \diamond G) = \sup_{\mathfrak{S}} Pr^{\mathfrak{S}}(s \models \diamond G).$$

where policy \mathfrak{S} ranges over all, infinitely (countably) many, policies.

Example

Equation system for max-reach probabilities

Equation system for max-reach probabilities

Let \mathcal{M} be a finite MDP with state space S , $s \in S$ and $G \subseteq S$. The vector $(x_s)_{s \in S}$ with $x_s = Pr^{\max}(s \models \diamond G)$ yields the unique solution of the following equation system:

- ▶ If $s \in G$, then $x_s = 1$.
- ▶ If $s \not\models \exists \diamond G$, then $x_s = 0$.
- ▶ If $s \models \exists \diamond G$ and $s \notin G$, then

$$x_s = \max \left\{ \sum_{t \in S} \mathbf{P}(s, \alpha, t) \cdot x_t \mid \alpha \in Act(s) \right\}$$

This is an instance of the Bellman equation for dynamic programming.

Maximal reachability probabilities

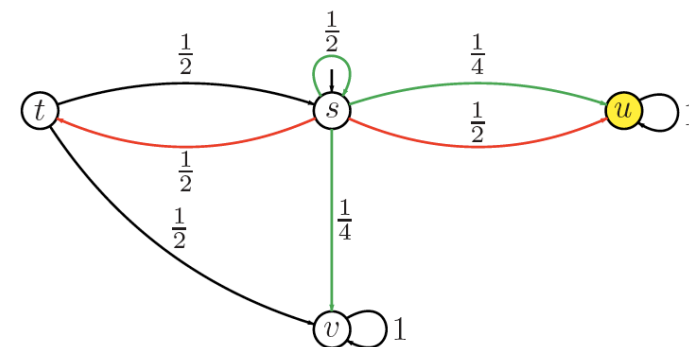
Minimal guarantees for safety properties

Reasoning about the maximal probabilities for $\diamond G$ is needed, e.g., for showing that $Pr^{\mathfrak{G}}(s \models \diamond G) \leq \varepsilon$ for all policies \mathfrak{G} and some small upper bound $0 < \varepsilon \leq 1$. Then:

$$Pr^{\mathfrak{G}}(s \models \square \neg G) \geq 1 - \varepsilon \quad \text{for all policies } \mathfrak{G}.$$

The task to compute $Pr^{\max}(s \models \diamond G)$ can thus be understood as showing that a safety property (namely $\square \neg G$) holds with sufficiently large probability, viz. $1 - \varepsilon$, regardless of the resolution of nondeterminism.

Example



equation system for reachability objective $\diamond \{ u \}$ is:

$$x_u = 1 \text{ and } x_v = 0$$

$$x_s = \max \left\{ \frac{1}{2}x_s + \frac{1}{4}x_u + \frac{1}{4}x_v, \frac{1}{2}x_u + \frac{1}{2}x_t \right\} \quad \text{and} \quad x_t = \frac{1}{2}x_s + \frac{1}{2}x_v$$

Value iteration

The previous theorem suggests to calculate the values

$$x_s = Pr^{\max}(s \models \diamond G)$$

by **successive approximation**.

For the states $s \models \exists \diamond G$ and $s \notin G$, we have $x_s = \lim_{n \rightarrow \infty} x_s^{(n)}$ where

$$x_s^{(0)} = 0 \quad \text{and} \quad x_s^{(n+1)} = \max \left\{ \sum_{t \in S} \mathbf{P}(s, \alpha, t) \cdot x_t^{(n)} \mid \alpha \in \text{Act}(s) \right\}.$$

Note that $x_s^{(0)} \leq x_s^{(1)} \leq x_s^{(2)} \leq \dots$. Thus, the values $Pr^{\max}(s \models \diamond G)$ can be approximated by successively computing the vectors

$$(x_s^{(0)}), (x_s^{(1)}), (x_s^{(2)}), \dots,$$

until $\max_{s \in S} |x_s^{(n+1)} - x_s^{(n)}|$ is below a certain (typically very small) threshold.

Equation system for min-reach probabilities

Equation system for min-reach probabilities

Let \mathcal{M} be a finite MDP with state space S , $s \in S$ and $G \subseteq S$. The vector $(x_s)_{s \in S}$ with $x_s = Pr^{\min}(s \models \diamond G)$ yields the unique solution of the following equation system:

- ▶ If $s \in G$, then $x_s = 1$.
- ▶ If $Pr^{\min}(s \models G) = 0$, then $x_s = 0$.
- ▶ If $Pr^{\min}(s \models G) > 0$ and $s \notin G$, then

$$x_s = \min \left\{ \sum_{t \in S} \mathbf{P}(s, \alpha, t) \cdot x_t \mid \alpha \in \text{Act}(s) \right\}$$

Positional policies suffice for reach probabilities

Existence of optimal positional policies

Let \mathcal{M} be a finite MDP with state space S , and $G \subseteq S$. There exists a **positional** policy \mathfrak{S} such that for any $s \in S$ it holds:

$$Pr^{\mathfrak{S}}(s \models \diamond G) = Pr^{\max}(s \models \diamond G).$$

Proof:

On the blackboard.

Preprocessing

The preprocessing required to compute the set

$$S_{=0}^{\min} = \{s \in S \mid Pr^{\min}(s \models \diamond G) = 0\}$$

can be performed by graph algorithms. The set $S_{=0}^{\min}$ is given by $S \setminus T$ where

$$T = \bigcup_{n \geq 0} T_n$$

and $T_0 = G$ and, for $n \geq 0$:

$$T_{n+1} = T_n \cup \{s \in S \mid \forall \alpha \in \text{Act}(s) \exists t \in T_n. \mathbf{P}(s, \alpha, t) > 0\}.$$

As $T_0 \subseteq T_1 \subseteq T_2 \subseteq \dots \subseteq S$ and S is finite, the sequence $(T_n)_{n \geq 0}$ eventually stabilizes, i.e., for some $n \geq 0$, $T_n = T_{n+1} = \dots = T$.

It follows: $Pr^{\min}(s \models \diamond G) > 0$ if and only if $s \in T$.

Preprocessing

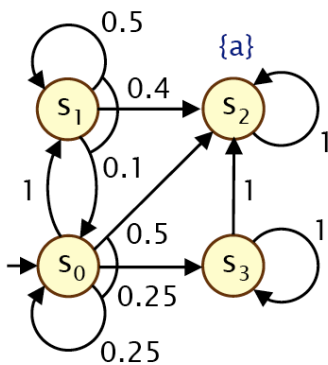
Algorithm 46 Computing the set of states s with $Pr^{\min}(s \models \Diamond B) = 0$

Input: finite MDP \mathcal{M} with state space S and $B \subseteq S$
Output: $\{s \in S \mid Pr^{\min}(s \models \Diamond B) = 0\}$

```

T := B;
R := B;
while R ≠ ∅ do
  let t ∈ R;
  R := R \ {t};
  for all (s, α) ∈ Pre(t) with s ∉ T do
    remove α from Act(s)
  if Act(s) = ∅ then
    add s to R and T
  fi
od
od
return T
    
```

Example value iteration



Determine $Pr^{\min}(s_i \models \Diamond \{s_2\})$.

Positional policies for min-reach probabilities

Existence of optimal positional policies

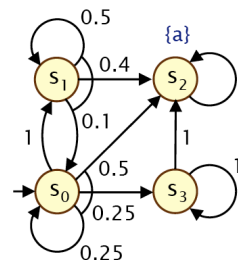
Let \mathcal{M} be a finite MDP with state space S , and $G \subseteq S$. There exists a **positional** policy \mathcal{G} such that for any $s \in S$ it holds:

$$Pr^{\mathcal{G}}(s \models \Diamond G) = Pr^{\min}(s \models \Diamond G).$$

Proof:

Similar to the case for maximal reachability probabilities.

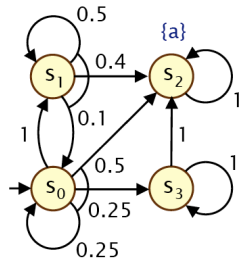
Example value iteration



Determine $Pr^{\min}(s_i \models \Diamond \{s_2\})$

1. $G = \{s_2\}, S_{=0}^{\min} = \{s_3\}, S \setminus (G \cup S_{=0}^{\min}) = \{s_0, s_1\}$.
2. $(x_s^{(0)}) = (0, 0, 1, 0)$
3. $(x_s^{(1)}) = (\min(1 \cdot 0, 0.25 \cdot 0 + 0.25 \cdot 0 + 0.5 \cdot 1),$
 $0.1 \cdot 0 + 0.5 \cdot 0 + 0.4 \cdot 1, 1, 0)$
4. $= (0, 0.4, 1, 0)$
5. $(x_s^{(2)}) = (\min(1 \cdot 0.4, 0.25 \cdot 0 + 0.25 \cdot 0 + 0.5 \cdot 1),$
 $0.1 \cdot 0 + 0.5 \cdot 0.4 + 0.4 \cdot 1, 1, 0)$
6. $= (0.4, 0.6, 1, 0)$
7. $(x_s^{(3)}) = \dots$

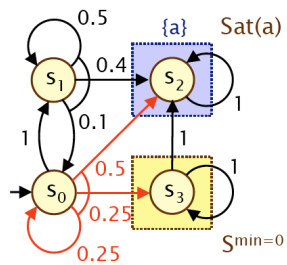
Example value iteration



Determine $Pr^{\min}(s_i \models \diamond\{s_2\})$

	$[x_0^{(n)}, x_1^{(n)}, x_2^{(n)}, x_3^{(n)}]$
n=0:	[0.000000, 0.000000, 1, 0]
n=1:	[0.000000, 0.400000, 1, 0]
n=2:	[0.400000, 0.600000, 1, 0]
n=3:	[0.600000, 0.740000, 1, 0]
n=4:	[0.650000, 0.830000, 1, 0]
n=5:	[0.662500, 0.880000, 1, 0]
n=6:	[0.665625, 0.906250, 1, 0]
n=7:	[0.666406, 0.919688, 1, 0]
n=8:	[0.666602, 0.926484, 1, 0]
...	
n=20:	[0.666667, 0.933332, 1, 0]
n=21:	[0.666667, 0.933332, 1, 0]
	$\approx [2/3, 14/15, 1, 0]$

Optimal positional policy



- ▶ Outcome of the value iteration $(x_s) = (\frac{2}{3}, \frac{14}{15}, 1, 0)$
- ▶ How to obtain the optimal policy from this result?
- ▶ $x_{s_0} = \min(1 \cdot \frac{14}{15}, 0.5 \cdot 1 + 0.25 \cdot 0 + 0.25 \cdot \frac{2}{3})$
 $\min(\frac{14}{15}, \frac{2}{3})$
- ▶ Thus the optimal policy always selects **red** in s_0 .

Optimal positional policy

Positional policies \mathcal{G}_{\min} and \mathcal{G}_{\max} thus yield:

$$Pr^{\mathcal{G}_{\min}}(s \models \diamond G) = Pr^{\min}(s \models \diamond G) \text{ for all states } s \in S$$

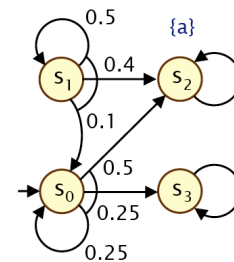
$$Pr^{\mathcal{G}_{\max}}(s \models \diamond G) = Pr^{\max}(s \models \diamond G) \text{ for all states } s \in S$$

These policies are obtained as follows:

$$\mathcal{G}_{\min}(s) = \arg \min \{ \sum_{t \in S} P(s, \alpha, t) \cdot Pr^{\min}(t \models \diamond G) \mid \alpha \in Act \}$$

$$\mathcal{G}_{\max}(s) = \arg \max \{ \sum_{t \in S} P(s, \alpha, t) \cdot Pr^{\max}(t \models \diamond G) \mid \alpha \in Act \}$$

Induced DTMC



- ▶ Outcome of the value iteration $(x_s) = (\frac{2}{3}, \frac{14}{15}, 1, 0)$
- ▶ How to obtain the optimal policy from this results?
- ▶ $x_{s_0} = \min(1 \cdot \frac{14}{15}, 0.5 \cdot 1 + 0.5 \cdot 0 + 0.25 \cdot \frac{2}{3})$
 $\min(\frac{14}{15}, \frac{2}{3})$
- ▶ Thus the optimal policy always selects **red**.

An alternative approach

A viable alternative to value iteration is **linear programming**.

Maximal reach probabilities as a linear program

Linear program for max-reach probabilities

Consider a finite MDP with state space S , and $G \subseteq S$. The values $x_s = Pr^{\max}(s \models \diamond G)$ are the unique solution of the **linear program**:

- ▶ If $s \in G$, then $x_s = 1$.
- ▶ If $s \not\models \exists \diamond G$, then $x_s = 0$.
- ▶ If $s \not\models \exists \diamond G$ and $s \notin G$, then $0 \leq x_s \leq 1$ and for all $\alpha \in Act(s)$:

$$x_s \geq \sum_{t \in S} P(s, \alpha, t) \cdot x_t$$

where $\sum_{s \in S} x_s$ is **minimal**.

Proof:

See lecture notes.

Linear programming

Linear programming

Optimisation of a linear objective function subject to linear (in)equalities.

Let x_1, \dots, x_n be real-valued variables. Maximise (or minimise) the **objective** function:

$$c_1 \cdot x_1 + c_2 \cdot x_2 + \dots + c_n \cdot x_n \quad \text{for constants } c_1, \dots, c_n \in \mathbb{R}$$

subject to the constraints

$$a_{11} \cdot x_1 + a_{12} \cdot x_2 + \dots + a_{1n} \cdot x_n \leq b_1$$

.....

$$a_{m1} \cdot x_1 + a_{m2} \cdot x_2 + \dots + a_{mn} \cdot x_n \leq b_m.$$

Solution techniques: e.g., Simplex, ellipsoid method, interior point method.

Minimal reach probabilities as a linear program

Linear program for min-reach probabilities

Consider a finite MDP with state space S , and $G \subseteq S$. The values $x_s = Pr^{\min}(s \models \diamond G)$ are the unique solution of the **linear program**:

- ▶ If $s \in G$, then $x_s = 1$.
- ▶ If $Pr^{\min}(s \models \diamond G) = 0$, then $x_s = 0$.
- ▶ If $Pr^{\min}(s \models \diamond G) > 0$ and $s \notin G$ then $0 \leq x_s \leq 1$ and for all $\alpha \in Act(s)$:

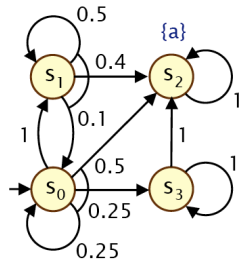
$$x_s \leq \sum_{t \in S} P(s, \alpha, t) \cdot x_t$$

where $\sum_{s \in S} x_s$ is **maximal**.

Proof:

See lecture notes.

Example linear programming

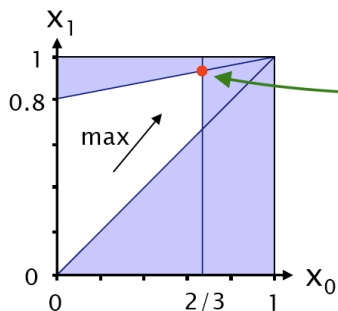
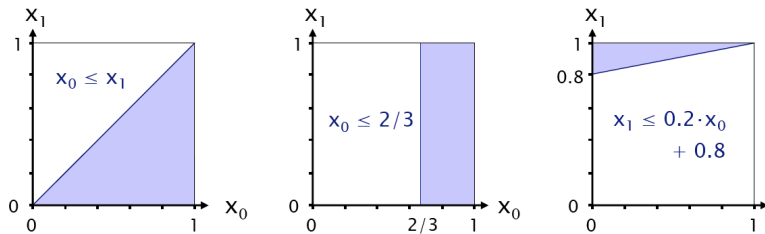


- $G = \{s_2\}, S_{=0}^{\min} = \{s_3\}, S \setminus (G \cup S_{=0}^{\min}) = \{s_0, s_1\}$.
- Maximise $x_0 + x_1$ subject to the constraints:

$$\begin{aligned} x_0 &\leq x_1 \\ x_0 &\leq \frac{1}{4} \cdot x_0 + \frac{1}{2} \\ x_1 &\leq \frac{1}{10} \cdot x_0 + \frac{1}{2} \cdot x_1 + \frac{2}{5} \end{aligned}$$

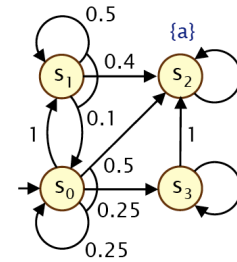
Determine $Pr^{\min}(s_i \models \diamond\{s_2\})$

Example linear programming



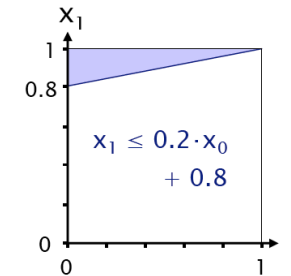
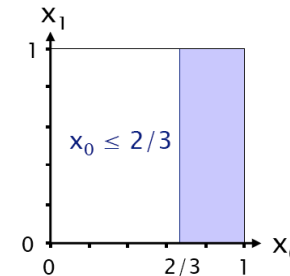
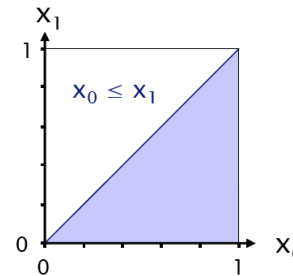
Solution:
 (x_0, x_1)
 $=$
 $(2/3, 14/15)$

Example linear programming

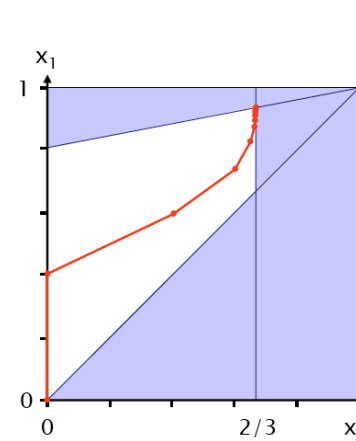


- $G = \{s_2\}, S_{=0}^{\min} = \{s_3\}, S \setminus (G \cup S_{=0}^{\min}) = \{s_0, s_1\}$.
- Maximise $x_0 + x_1$ subject to the constraints:

$$\begin{aligned} x_0 &\leq x_1 \\ x_0 &\leq \frac{2}{3} \\ x_1 &\leq \frac{2}{5} \cdot x_0 + \frac{4}{5} \end{aligned}$$



Value iteration vs. linear programming



	$[x_0^{(n)}, x_1^{(n)}, x_2^{(n)}, x_3^{(n)}]$
n=0:	[0.000000, 0.000000, 1, 0]
n=1:	[0.000000, 0.400000, 1, 0]
n=2:	[0.400000, 0.600000, 1, 0]
n=3:	[0.600000, 0.740000, 1, 0]
n=4:	[0.650000, 0.830000, 1, 0]
n=5:	[0.662500, 0.880000, 1, 0]
n=6:	[0.665625, 0.906250, 1, 0]
n=7:	[0.666406, 0.919688, 1, 0]
n=8:	[0.666602, 0.926484, 1, 0]
...	
n=20:	[0.666667, 0.933332, 1, 0]
n=21:	[0.666667, 0.933332, 1, 0]

$\approx [2/3, 14/15, 1, 0]$

This curve shows how the value iteration approach approximates the solution.

Time complexity

Time complexity

For finite MDP \mathcal{M} with state space S , $G \subseteq S$ and $s \in S$, the values $Pr^{\max}(s \models \diamond G)$ can be computed in time polynomial in the size of \mathcal{M} .
The same holds for $Pr^{\min}(s \models \diamond G)$.

Proof:

Thanks to the characterisation as a linear program and polynomial time techniques to solve such linear programs such as ellipsoid methods.

Corollary

For finite MDPs, the question whether $Pr^{\mathfrak{G}}(s \models \diamond G) \leq p$ for some rational $p \in [0, 1[$ is decidable in polynomial time.

Policy iteration

Value iteration

In value iteration, we iteratively attempt to improve the minimal (or maximal) reachability probabilities by starting with an underestimation, viz. zero for all states.

Policy iteration

In **policy** iteration, the idea is to start with an arbitrary positional policy and improve it in a step-by-step fashion, so as to determine the optimal one.

Yet another alternative approach

A viable alternative to value iteration and linear programming is **policy iteration**.

Policy iteration

Policy iteration

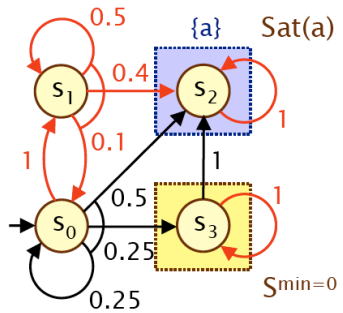
1. Start with an arbitrary positional policy \mathfrak{G} that selects some $\alpha \in Act(s)$ for each state s .
2. Compute the reachability probabilities $Pr^{\mathfrak{G}}(s \models \diamond G)$. This amounts to solving a linear equation system on DTMC $\mathcal{M}_{\mathfrak{G}}$.
3. Improve the policy in every state according to the following rules:

$$\mathfrak{G}^{(i+1)}(s) = \arg \min \left\{ \sum_{t \in S} \mathbf{P}(s, \alpha, t) \cdot Pr^{\mathfrak{G}^{(i)}}(t \models \diamond G) \mid \alpha \in Act \right\} \text{ or}$$

$$\mathfrak{G}^{(i+1)}(s) = \arg \max \left\{ \sum_{t \in S} \mathbf{P}(s, \alpha, t) \cdot Pr^{\mathfrak{G}^{(i)}}(t \models \diamond G) \mid \alpha \in Act \right\}$$

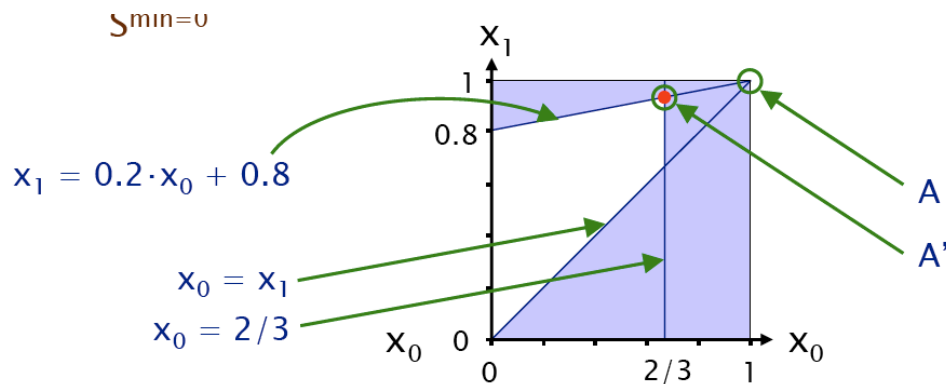
4. Repeat steps 2. and 3. until the policy does not change.

Policy iteration: example



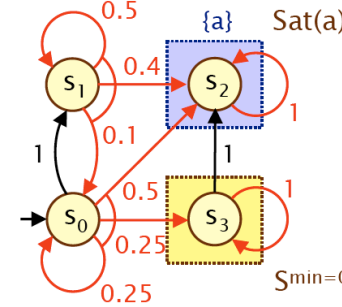
- ▶ Let $G = \{s_2\}$.
- ▶ Consider an arbitrary policy \mathfrak{G} .
- ▶ Compute $x_i = Pr^{\mathfrak{G}}(s_i \models \Diamond G)$ for all i .
- ▶ Then: $x_2 = 1, x_3 = 0$,
and $x_0 = x_1, x_1 = \frac{1}{10} \cdot x_0 + \frac{1}{2} \cdot x_1 + \frac{2}{5}$.
- ▶ This yields $x_0 = x_1 = x_2 = 1$ and $x_3 = 0$.
- ▶ Change policy \mathfrak{G} in s_0 , yielding policy \mathfrak{G}' .
- ▶ This yields $\min(1 \cdot 1, \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 0 + \frac{1}{4} \cdot 1)$
that is, $\min(1, \frac{3}{4}) = \frac{3}{4}$.

Graphical representation of policy iteration



where A denotes policy \mathfrak{G} and A' policy \mathfrak{G}' .

Policy iteration: example



- ▶ Let $G = \{s_2\}$.
- ▶ Consider the adapted policy \mathfrak{G}' .
- ▶ Compute $x_i = Pr^{\mathfrak{G}'}(s_i \models \Diamond G)$ for all i .
- ▶ Then: $x_2 = 1, x_3 = 0$,
and $x_0 = \frac{1}{4} \cdot x_0 + \frac{1}{2}, x_1 = \frac{1}{10} \cdot x_0 + \frac{1}{2} \cdot x_1 + \frac{2}{5}$.
- ▶ This yields $x_0 = \frac{2}{3}, x_1 = \frac{14}{15}, x_2 = 1$ and $x_3 = 0$.
- ▶ This policy is optimal.

Overview

- 1 Markov Decision Processes
- 2 Probabilities in MDPs
- 3 Policies
 - Positional policies
 - Finite-memory policies
- 4 Reachability probabilities
 - Mathematical characterisation
 - Value iteration
 - Linear programming
 - Policy iteration
- 5 Summary

Summary

Important points

1. Maximal reachability probabilities are suprema over reachability probabilities for all, potentially infinitely many, policies.
2. They are characterised by equation systems with maximal operators.
3. There exists a positional policy that yields the maximal reachability probability.
4. Such policies can be determined using value or policy iteration.
5. Or, alternatively, in polynomial time using linear programming.
6. Positional policies are not powerful enough for arbitrary ω -regular properties.