

Planning and Optimization

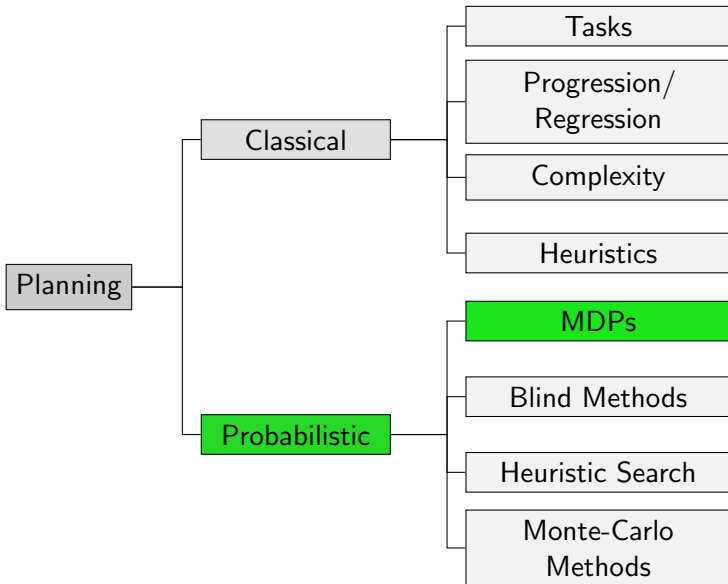
F2. Policies & Compact Description

Gabriele Röger and Thomas Keller

Universität Basel

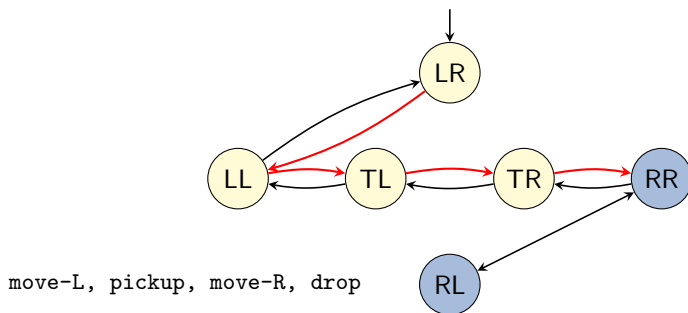
November 21, 2018

Content of this Course



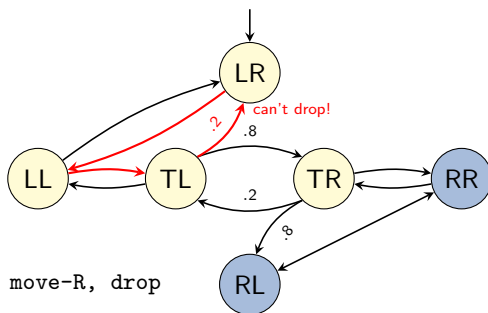
Policies & Value Functions

Solutions in SSPs



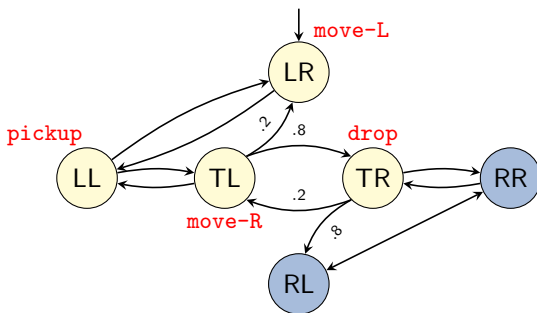
- solution in deterministic transition systems is **plan**, i.e., a goal path from s_0 to some $s_* \in S_*$
- **cheapest** plan is **optimal solution**
- deterministic agent that **executes** plan will reach goal

Solutions in SSPs



- probabilistic agent **will not reach goal** or **cannot execute** plan
- non-determinism can lead to **different outcome** than **anticipated** in plan
- require a more general solution: a **policy**

Solutions in SSPs



- policy must be allowed to be **cyclic**
- policy must be able to **branch** over outcomes
- policy assigns **applicable labels** to states

Policy for SSPs

Definition (Policy for SSPs)

Let $\mathcal{T} = \langle S, L, c, T, s_0, S_\star \rangle$ be an SSP. A **policy** for \mathcal{T} is a mapping $\pi : S \rightarrow L \cup \{\perp\}$ such that $\pi(s) \in L(s) \cup \{\perp\}$ for all s .

The set of **reachable states** $S_\pi(s)$ from s under π is defined recursively as the smallest set satisfying the rules

- $s \in S_\pi(s)$ and
- $\text{succ}(s', \pi(s')) \subseteq S_\pi(s)$ for all $s' \in S_\pi(s) \setminus S_\star$ where $\pi(s') \neq \perp$.

If $\pi(s') \neq \perp$ for all $s' \in S_\pi(s)$, then π is **executable in s** .

Policy Representation

- size of **explicit representation** of executable policy π is $|\mathcal{S}_\pi(s_0)|$
- often, $|\mathcal{S}_\pi(s_0)|$ similar to $|\mathcal{S}|$
- **compact** policy representation, e.g. via value function approximation or neural networks, is active research area
⇒ not covered in this course
- instead, we consider **small state spaces** for basic algorithms
- or **online** planning where planning for the current state s_0 is interleaved with **execution** of $\pi(s_0)$

Value Functions of SSPs

Definition (Value Functions of SSPs)

Let $\mathcal{T} = \langle S, L, c, T, s_0, S_\star \rangle$ be an SSP and π be an executable policy for \mathcal{T} . The **state-value** $V_\pi(s)$ of s under π is defined as

$$V_\pi(s) := \begin{cases} 0 & \text{if } s \in S_\star \\ Q_\pi(s, \pi(s)) & \text{otherwise,} \end{cases}$$

where the **action-value** $Q_\pi(s, \ell)$ under π is defined as

$$Q_\pi(s, \ell) := c(\ell) + \sum_{s' \in \text{succ}(s, \ell)} (T(s, \ell, s') \cdot V_\pi(s')).$$

Example: Value Functions of SSPs

Example

Consider example task and π with $\pi(\text{LR}) = \text{move-L}$,
 $\pi(\text{LL}) = \text{pickup}$, $\pi(\text{TL}) = \text{move-R}$ and $\pi(\text{TR}) = \text{drop}$.

$$V_*(\text{LR}) = 1 + V_*(\text{LL})$$

$$V_*(\text{LL}) = 1 + V_*(\text{TL})$$

$$V_*(\text{TL}) = 1 + (0.8 \cdot V_*(\text{RR})) + (0.2 \cdot V_*(\text{LR}))$$

$$V_*(\text{TR}) = 1 + V_*(\text{RR})$$

$$V_*(\text{RL}) = 0$$

$$V_*(\text{RR}) = 0$$

What is the solution of this? \Rightarrow next week!

Bellman Optimality Equation

Definition (Optimal Policy in SSPs)

Let the **Bellman optimality equation** for a state s of an SSP be the set of equations that describes $V_*(s)$, where

$$V_*(s) := \begin{cases} 0 & \text{if } s \in S_* \\ \min_{\ell \in L(s)} Q_*(s, \ell) & \text{otherwise,} \end{cases}$$

$$Q_*(s, \ell) := c(\ell) + \sum_{s' \in \text{succ}(s, \ell)} (T(s, \ell, s') \cdot V_*(s')).$$

A policy π^* is an **optimal policy** if $\pi^*(s) \in \arg \min_{\ell \in L(s)} Q_*(s, \ell)$ for all $s \in S$, and the **expected cost** of π^* in \mathcal{T} is $V_*(s_0)$.

Dead-end States

- **dead-end states** are a problem with our formalization
- each policy with **non-zero probability** of reaching a dead-end has **infinite state-value**
- one solution is to search for policy with **highest probability to reach the goal**
- unfortunately, this **ignores costs**
- there is also research on **dead-end detection**
- in this course, we only consider SSPs, FH-MDPs and DR-MDPs that are **dead-end free**

Policies for FH-MDPs

- What is the optimal policy for the SSP at the [blackboard](#)?

Policies for FH-MDPs

- What is the optimal policy for the SSP at the [blackboard](#)?
- Can we do better if we regard this as an FH-MDP?

Policies for FH-MDPs

- What is the optimal policy for the SSP at the **blackboard**?
- Can we do better if we regard this as an FH-MDP?
- Yes, by acting differently **close to the horizon**.

Policy for FH-MDPs

Definition (Policy for FH-MDPs)

Let $\mathcal{T} = \langle S, L, R, T, s_0, H \rangle$ be an FH-MDP. A policy for \mathcal{T} is a mapping $\pi : S \times \{1, \dots, H\} \rightarrow L \cup \{\perp\}$ such that $\pi(s, d) \in L(s) \cup \{\perp\}$ for all s .

The set of reachable states $S_\pi(s, d)$ from s with d steps-to-go under π is defined recursively as the smallest set satisfying the rules

- $\langle s, d \rangle \in S_\pi(s, d)$ and
- $\langle s'', d' - 1 \rangle \in S_\pi(s, d)$ for all $s'' \in \text{succ}(s', \pi(s'))$ and $\langle s', d' \rangle \in S_\pi(s)$ with $d' > 0$ and $\pi(s', d') \neq \perp$.

If $\pi(s', d') \neq \perp$ for all $\langle s', d' \rangle \in S_\pi(s, d)$ with $d' > 0$, then π is executable in s .

Value Functions for FH-MDPs

Definition (Value Functions for FH-MDPs)

Let $\mathcal{T} = \langle S, L, c, T, s_0, H \rangle$ be an FH-MDP and π be an executable policy for \mathcal{T} . The state-value $V_\pi(s, d)$ of s and d under π is defined as

$$V_\pi(s, d) := \begin{cases} 0 & \text{if } d = 0 \\ Q_\pi(s, d, \pi(s)) & \text{otherwise,} \end{cases}$$

where the action-value $Q_\pi(s, d, \ell)$ under π is defined as

$$Q_\pi(s, d, \ell) := R(s, \ell) + \sum_{s' \in \text{succ}(s, \ell)} (T(s, \ell, s') \cdot V_\pi(s', d - 1)).$$

Bellman Optimality Equation

Definition (Optimal Policy in FH-MDPs)

Let the Bellman optimality equation for a state s of an FH-MDP be the set of equations that describes $V_*(s, d)$, where

$$V_*(s, d) := \begin{cases} 0 & \text{if } d = 0 \\ \max_{\ell \in L(s)} Q_*(s, d, \ell) & \text{otherwise,} \end{cases}$$

$$Q_*(s, d, \ell) := R(s, \ell) + \sum_{s' \in \text{succ}(s, \ell)} (T(s, \ell, s') \cdot V_*(s', d - 1)).$$

A policy π^* is an optimal policy if

$\pi^*(s, d) \in \arg \max_{\ell \in L(s)} Q_*(s, d, \ell)$ for all $s \in S$ and $d \in \{1, \dots, H\}$, and the **expected reward** of π^* in \mathcal{T} is $V_*(s_0, H)$.

(Optimal) Policy and Value Functions for DR-MDPs

- policy does **not distinguish states** based on steps-to-go (or rather the reverse “distance-from-init”)
- value functions have no **“terminal case”**
- value functions **discount** reward with γ
- Bellman optimality equation derived from value functions **as for FH-MDP**

Factored MDPs

Factored SSPs

We would like to specify huge SSPs without enumerating states. In classical planning, we achieved this via propositional planning tasks:

- represent different aspects of the world
in terms of different **Boolean state variables**
- treat state variables as atomic propositions
↪ a state is a **valuation of state variables**
- n state variables induce 2^n states
↪ **exponentially more compact** than “flat” representations

⇒ can also be used for SSPs

Reminder: Syntax of Operators

Definition (Operator)

An **operator** o over state variables V is an object with three properties:

- a **precondition** $pre(o)$, a logical formula over V
- an **effect** $eff(o)$ over V , defined on the following slides
- a **cost** $cost(o) \in \mathbb{R}_0^+$

⇒ can also be used for SSPs

Reminder: Syntax of Effects

Definition (Effect)

Effects over state variables V are inductively defined as follows:

- If $v \in V$ is a state variable, then v and $\neg v$ are effects (**atomic effect**).
- If e_1, \dots, e_n are effects, then $(e_1 \wedge \dots \wedge e_n)$ is an effect (**conjunctive effect**).

The special case with $n = 0$ is the **empty effect** \top .

- If χ is a logical formula and e is an effect, then $(\chi \triangleright e)$ is an effect (**conditional effect**).

Parentheses can be omitted when this does not cause ambiguity.

Syntax of Probabilistic Effects

Definition (Effect)

Effects over state variables V are inductively defined as follows:

- If $v \in V$ is a state variable, then v and $\neg v$ are effects (**atomic effect**).
- If e_1, \dots, e_n are effects, then $(e_1 \wedge \dots \wedge e_n)$ is an effect (**conjunctive effect**).

The special case with $n = 0$ is the **empty effect** \top .

- If χ is a logical formula and e is an effect, then $(\chi \triangleright e)$ is an effect (**conditional effect**).
- If e_1, \dots, e_n are effects and $p_1, \dots, p_n \in [0, 1]$ such that $\sum_{i=1}^n p_i = 1$, then $(p_1 : e_1 \mid \dots \mid p_n : e_n)$ is an effect (**probabilistic effect**).

Parentheses can be omitted when this does not cause ambiguity.

- **FDR tasks** can be generalized to SSPs in the same way
- both propositional and FDR tasks can be **generalized to FH-MDPs and DR-MDPs**

Summary

Summary

- Policies consider branching and cycles
- State-value of a policy describes **expected reward** of following that policy
- Related **Bellman optimality equation** describes optimal policy
- **Compact descriptions** that induce SSPs and MDPs analogous to classical planning