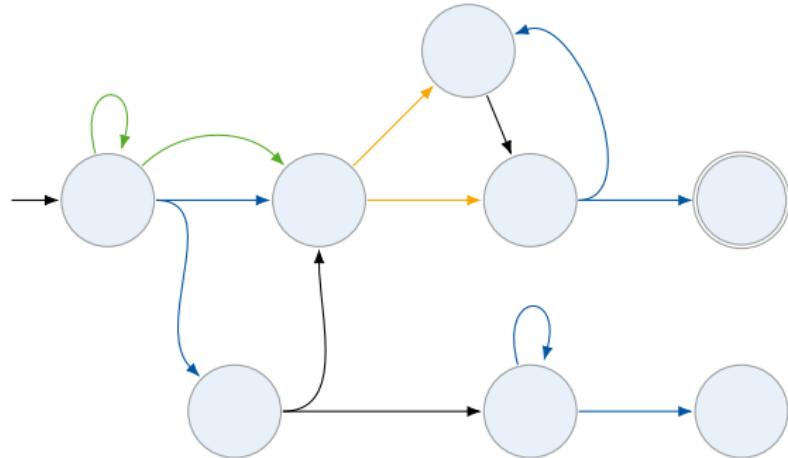


Learning Generalized Policies for Fully Observable Non-Deterministic Planning Domains

IJCAI 2024

Till Hofmann, Hector Geffner

August 6, 2024

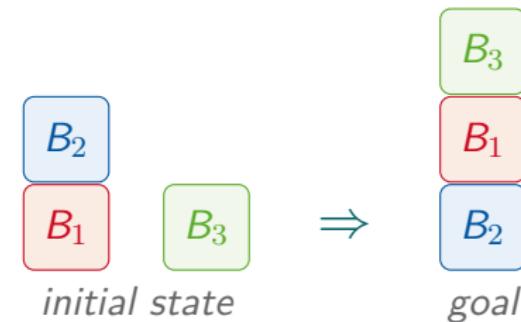


Generalized Policies for Classical Planning

Classical planning:

- ▶ Given: action model + initial state + goal
- ▶ Task: find a (concrete) policy $\pi_P : S \rightarrow 2^A$ that maps each state s to a set of actions $\pi_P(s)$

```
(:action pickup :parameters (?x)
  :precondition (and
    (clear ?x) (ontable ?x) (handempty))
  :effect (and (not (ontable ?x))
    (not (clear ?x)) (not (handempty))
    (holding ?x)))
(:action putdown ...)
(:action stack ...)
(:action unstack ...)
```



Generalized Policies for Classical Planning

Classical planning:

- ▶ Given: action model + initial state + goal
- ▶ Task: find a (concrete) policy $\pi_P : S \rightarrow 2^A$ that maps each state s to a set of actions $\pi_P(s)$

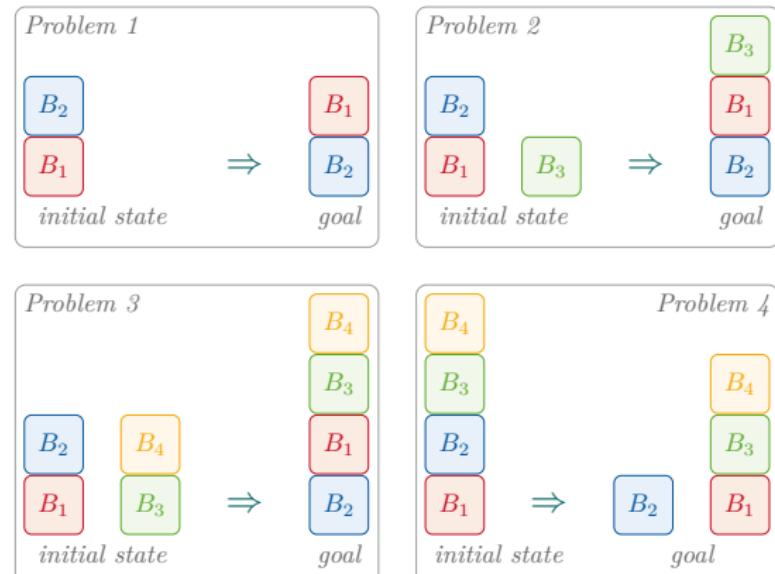
Generalized planning:

- ▶ Given: action model + **class of problems** \mathcal{Q}
- ▶ *Each* problem: initial state + goal
- ▶ Task: find a general policy that maps each problem $P \in \mathcal{Q}$ to a concrete policy $\pi(P) = \pi_P$

Example (Blocks World)

General policy:

- ① Unstack all misplaced blocks
- ② Stack them in the correct order



Generalized Policies for Classical Planning

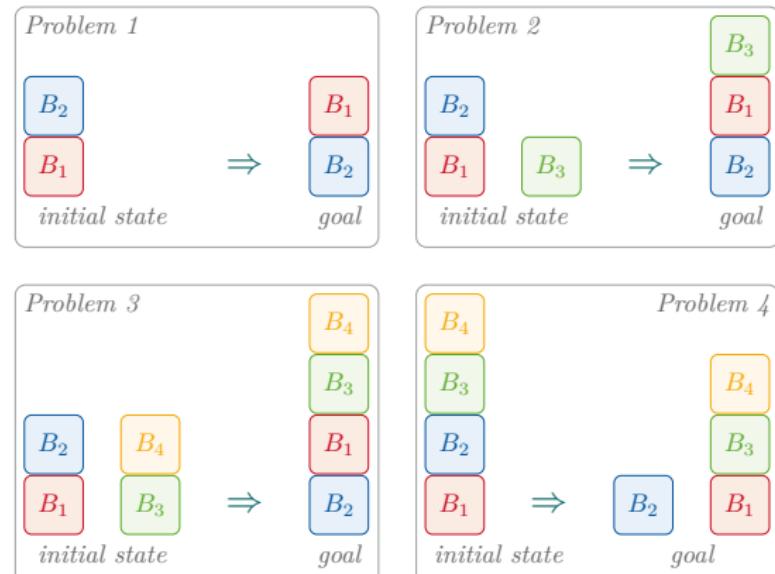
Classical planning:

- ▶ Given: action model + initial state + goal
- ▶ Task: find a (concrete) policy $\pi_P : S \rightarrow 2^A$ that maps each state s to a set of actions $\pi_P(s)$

Generalized planning:

- ▶ Given: action model + **class of problems** \mathcal{Q}
- ▶ *Each* problem: initial state + goal
- ▶ Task: find a general policy that maps each problem $P \in \mathcal{Q}$ to a concrete policy $\pi(P) = \pi_P$

assumed
to be closed



Example (Blocks World)

General policy:

- ① Unstack all misplaced blocks
- ② Stack them in the correct order

Expressing and Learning Generalized Policies

Rule-based general policies (Bonet and Geffner 2018):

- ▶ Express a policy as a set of rules of the form $C \mapsto E_1 \mid \dots \mid E_k$ over a feature space Φ
- ▶ Features: Boolean and numerical features, e.g., $\# \text{ clear blocks}$
- ▶ Learnable from a set of small problems (Bonet et al. 2019; Francès et al. 2021)

Example

General policy for clearing a one block x :

$$r_1 : \{\neg H, n > 0\} \mapsto \{H, n \downarrow\}$$

$$r_2 : \quad \quad \quad \{H\} \mapsto \{\neg H\}$$

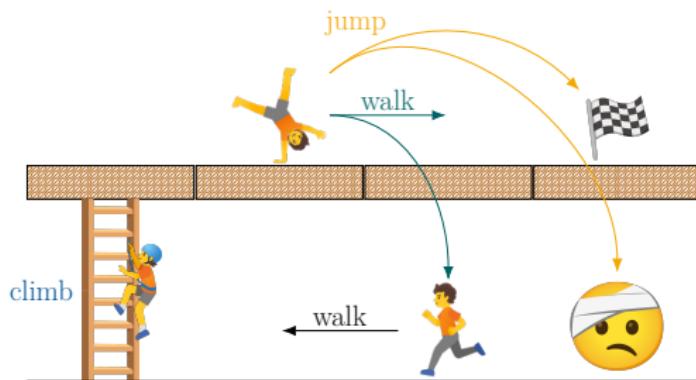
where

$H \triangleq$ true if holding a block

$n \triangleq$ number of blocks on top of x

Generalized Policies for FOND Planning

- ▶ FOND: action a may have multiple outcomes a_1, a_2, \dots not under the agent's control
- ▶ Fairness: for an infinitely occurring action a , all outcomes a_1, a_2, \dots occur infinitely often
- ▶ Strong-cyclic: all max. fair trajectories reach goal

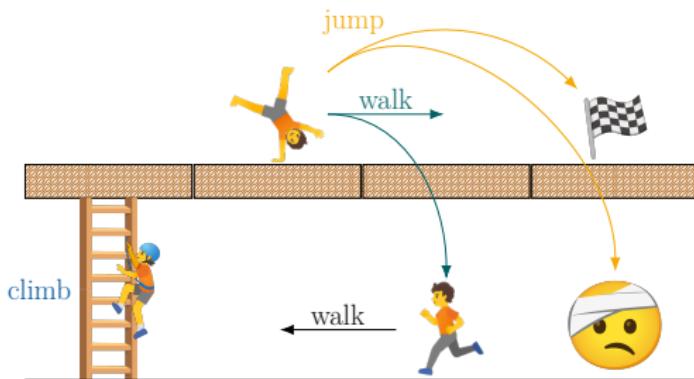


Generalized Policies for FOND Planning

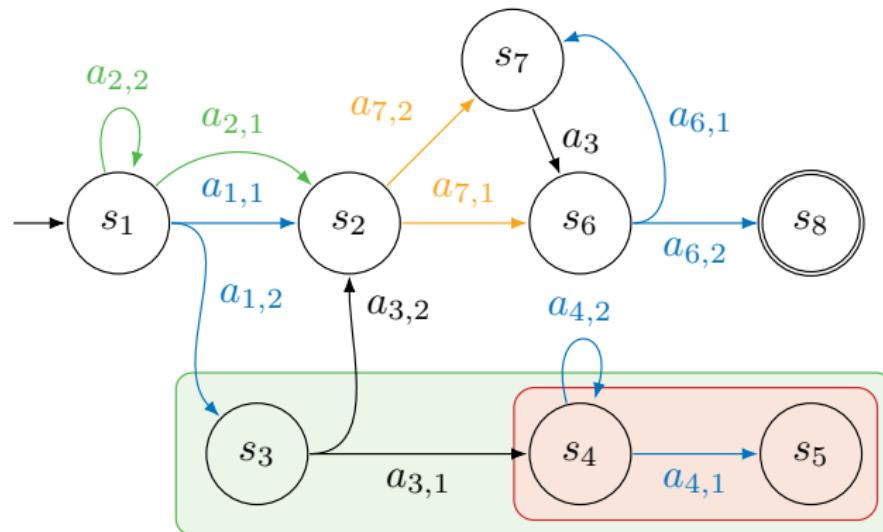
- ▶ FOND: action a may have multiple outcomes a_1, a_2, \dots not under the agent's control
- ▶ Fairness: for an infinitely occurring action a , all outcomes a_1, a_2, \dots occur infinitely often
- ▶ Strong-cyclic: all max. fair trajectories reach goal

Questions

- ① How do we express general policies for FOND domains?
- ② How can we learn FOND general policies?

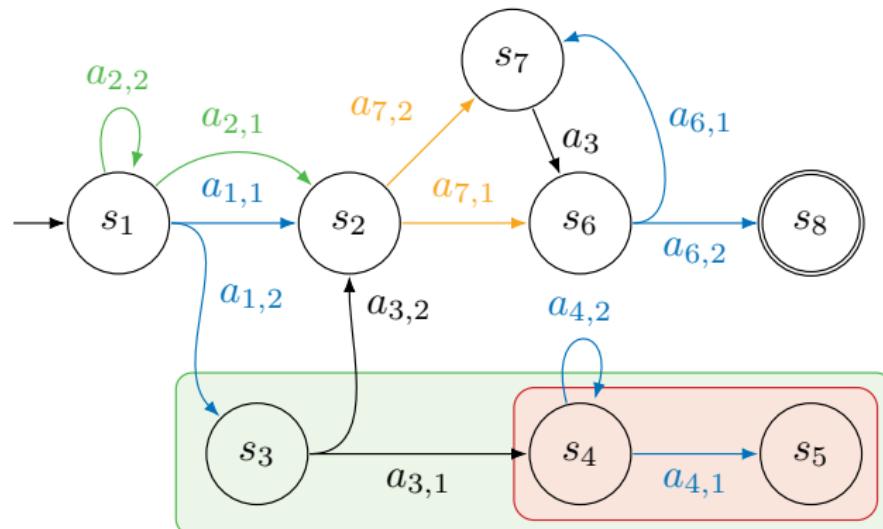


From Classical to FOND General Policies



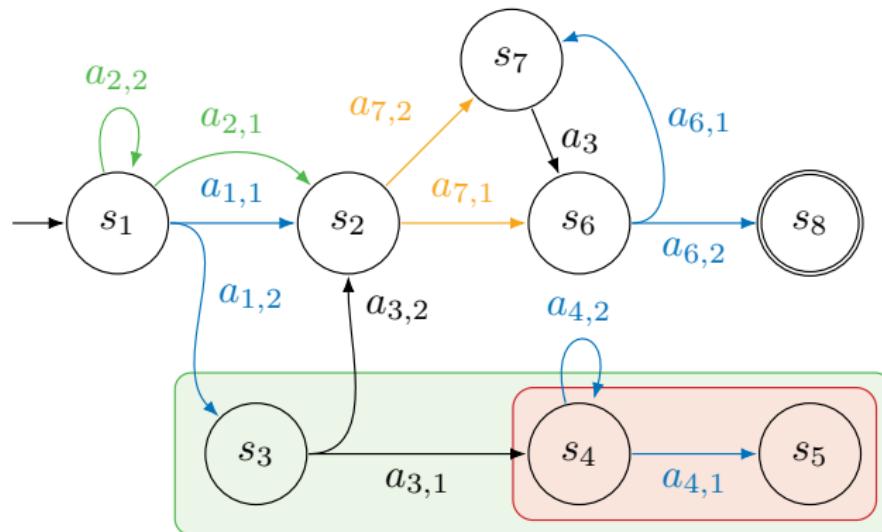
From Classical to FOND General Policies

All-outcome relaxation: agent can choose action outcome \rightarrow classical planning



FOND dead-ends relaxation dead-ends

From Classical to FOND General Policies

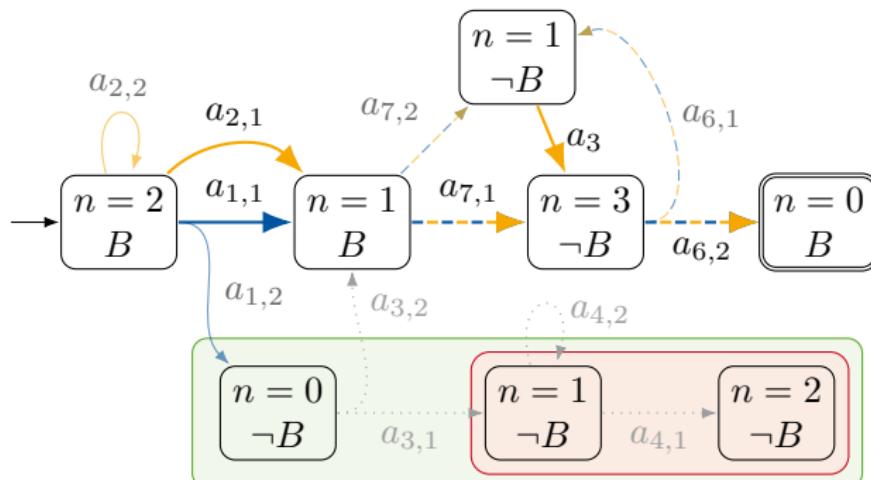


All-outcome relaxation: agent can choose action outcome \rightarrow classical planning

Observation

A general policy for the *all-outcome relaxation* is also a general policy for the original FOND problem if it avoids FOND dead-ends

From Classical to FOND General Policies



→ Policy 1

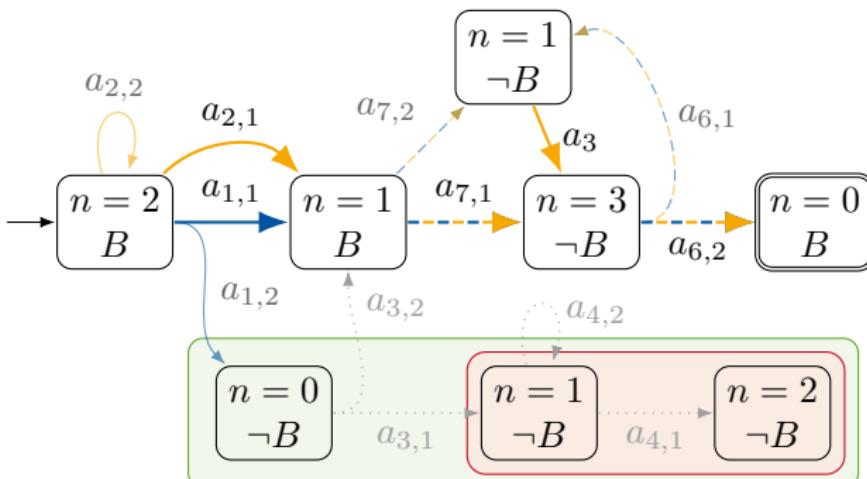
→ Policy 2

All-outcome relaxation: agent can choose action outcome → classical planning

Observation

A general policy for the *all-outcome relaxation* is also a general policy for the original FOND problem if it avoids FOND dead-ends

From Classical to FOND General Policies



→ Policy 1

→ Policy 2

All-outcome relaxation: agent can choose action outcome → classical planning

Observation

A general policy for the *all-outcome relaxation* is also a general policy for the original FOND problem if it avoids FOND dead-ends

Learning:

- ▶ Select features that describe the FOND dead-ends
- ▶ Learn a general policy for the relaxation that avoids dead-ends
- ⇒ Combinatorial optimization problem

more details on poster

Example: Learned Policy for Acrobatics

Selected features:

- ① Distance $d \equiv \text{dist}(\text{position}, \text{next-fwd}, \text{position}_G)$ between the current position and the goal position
- ② Boolean feature $U \equiv |\text{up}|$ which is true if the agent is currently on the beam
- ③ Boolean feature $B \equiv |\text{broken-leg}|$ which is true if the agent's leg is broken

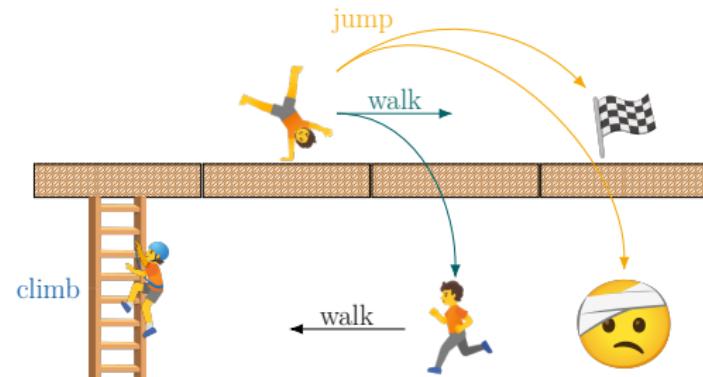
Learned policy rules:

$$r_1 : \{U, d > 0, \neg B\} \mapsto \{d \downarrow\}$$

$$r_2 : \{\neg B, \neg U\} \mapsto \{U\} \mid \{d \uparrow\}$$

Learned state constraint:

$$b_1 : \{B, \neg U\}$$



Evaluation

\mathcal{Q}	# problems	# solved	# training p.	# obj (train)	# obj (all)	t_{solve}/s	# features	# sel. features	# constraints	max complex.	policy cost
acrobatics	18	18	3	3	9	0.12	23	3	1	4	6
beam-walk	9	9	2	3	9	0.05	22	2	0	4	5
blocks3ops	95	95	4	4	20	3:44	194	3	0	5	11
blocks-clear	95	95	2	3	20	1.5	34	2	0	4	6
blocks-on	190	190	2	3	20	1:56	704	3	0	6	11
doors	19	19	5	7	33	1:18	625	4	1	10	19
islands	300	300	4	32	83	1h	1182	4	1	7	13
first-resp	99	15^M	2	5	36	33:40	332	5	2	7	20
miner	69	13^I	2	9	184	297h	1073	8	4	6	28
spiky-tire	170	36^I	3	6	23	1h	479	8	5	8	36
tireworld	980	7^C	1	3	100	0.11	27	5	4	4	12
triangle-tire	10	1^I	1	6	231	0.29	27	3	1	4	9

Evaluation

\mathcal{Q}	# problems	# solved	# training p.	# obj (train)	# obj (all)	$t_{\text{solve}}/\text{s}$	# features	# sel. features	# constraints	max complex.	policy cost
acrobatics	18	18	3	3	9	0.12	23	3	1	4	6
beam-walk	9	9	2	3	9	0.05	22	2	0	4	5
blocks3ops	95	95	4	4	20	3:44	194	3	0	5	11
blocks-clear	95	95	2	3	20	1.5	34	2	0	4	6
blocks-on	190	190	2	3	20	1:56	704	3	0	6	11
doors	19	19	5	7	33	1:18	625	4	1	10	19
islands	300	300	4	32	83	1h	1182	4	1	7	13
first-resp	99	15^M	2	5	36	33:40	332	5	2	7	20
miner	69	13^I	2	9	184	297h	1073	8	4	6	28
spiky-tire	170	36^I	3	6	23	1h	479	8	5	8	36
tireworld	980	7^C	1	3	100	0.11	27	5	4	4	12
triangle-tire	10	1^I	1	6	231	0.29	27	3	1	4	9

Evaluation

\mathcal{Q}	# problems	# solved	# training p.	# obj (train)	# obj (all)	t_{solve}/s	# features	# sel. features	# constraints	max complex.	policy cost
acrobatics	18	18	3	3	9	0.12	23	3	1	4	6
beam-walk	9	9	2	3	9	0.05	22	2	0	4	5
blocks3ops	95	95	4	4	20	3:44	194	3	0	5	11
blocks-clear	95	95	2	3	20	1.5	34	2	0	4	6
blocks-on	190	190	2	3	20	1:56	704	3	0	6	11
doors	19	19	5	7	33	1:18	625	4	1	10	19
islands	300	300	4	32	83	1h	1182	4	1	7	13
first-resp	99	15^M	2	5	36	33:40	332	5	2	7	20
miner	69	13^I	2	9	184	297h	1073	8	4	6	28
spiky-tire	170	36^I	3	6	23	1h	479	8	5	8	36
tireworld	980	7^C	1	3	100	0.11	27	5	4	4	12
triangle-tire	10	1^I	1	6	231	0.29	27	3	1	4	9

Evaluation

\mathcal{Q}	# problems	# solved	# training p.	# obj (train)	# obj (all)	t_{solve}/s	# features	# sel. features	# constraints	max complex.	policy cost
acrobatics	18	18	3	3	9	0.12	23	3	1	4	6
beam-walk	9	9	2	3	9	0.05	22	2	0	4	5
blocks3ops	95	95	4	4	20	3:44	194	3	0	5	11
blocks-clear	95	95	2	3	20	1.5	34	2	0	4	6
blocks-on	190	190	2	3	20	1:56	704	3	0	6	11
doors	19	19	5	7	33	1:18	625	4	1	10	19
islands	300	300	4	32	83	1h	1182	4	1	7	13
first-resp	99	15^M	2	5	36	33:40	332	5	2	7	20
miner	69	13^I	2	9	184	297h	1073	8	4	6	28
spiky-tire	170	36^I	3	6	23	1h	479	8	5	8	36
tireworld	980	7^C	1	3	100	0.11	27	5	4	4	12
triangle-tire	10	1^I	1	6	231	0.29	27	3	1	4	9

Evaluation

\mathcal{Q}	# problems	# solved	# training p.	# obj (train)	# obj (all)	t_{solve}/s	# features	# sel. features	# constraints	max complex.	policy cost
acrobatics	18	18	3	3	9	0.12	23	3	1	4	6
beam-walk	9	9	2	3	9	0.05	22	2	0	4	5
blocks3ops	95	95	4	4	20	3:44	194	3	0	5	11
blocks-clear	95	95	2	3	20	1.5	34	2	0	4	6
blocks-on	190	190	2	3	20	1:56	704	3	0	6	11
doors	19	19	5	7	33	1:18	625	4	1	10	19
islands	300	300	4	32	83	1h	1182	4	1	7	13
first-resp	99	15^M	2	5	36	33:40	332	5	2	7	20
miner	69	13^I	2	9	184	297h	1073	8	4	6	28
spiky-tire	170	36^I	3	6	23	1h	479	8	5	8	36
tireworld	980	7^C	1	3	100	0.11	27	5	4	4	12
triangle-tire	10	1^I	1	6	231	0.29	27	3	1	4	9

Evaluation

\mathcal{Q}	# problems	# solved	# training p.	# obj (train)	# obj (all)	$t_{\text{solve}}/\mathcal{S}$	# features	# sel. features	# constraints	max complex.	policy cost
proven to generalize to all instances											
acrobatics	18	18	3	3	9	0.12	23	3	1	4	6
beam-walk	9	9	2	3	9	0.05	22	2	0	4	5
blocks3ops	95	95	4	4	20	3:44	194	3	0	5	11
blocks-clear	95	95	2	3	20	1.5	34	2	0	4	6
blocks-on	190	190	2	3	20	1:56	704	3	0	6	11
doors	19	19	5	7	33	1:18	625	4	1	10	19
islands	300	300	4	32	83	1h	1182	4	1	7	13
first-resp	99	15^M	2	5	36	33:40	332	5	2	7	20
miner	69	13^I	2	9	184	297h	1073	8	4	6	28
spiky-tire	170	36^I	3	6	23	1h	479	8	5	8	36
tireworld	980	7^C	1	3	100	0.11	27	5	4	4	12
triangle-tire	10	1^I	1	6	231	0.29	27	3	1	4	9

LEARNING GENERALIZED POLICIES FOR FULLY OBSERVABLE NON-DETERMINISTIC PLANNING DOMAINS

Till Hofmann and Hector Geffner
RWTH Aachen University



Generalized FOND Planning

- Actions may have multiple (non-deterministic) outcomes, chosen by the environment
- Fairness: an infinity occurring action a , all outcomes a' must be equally likely
- Shortest path: all max. for implications road map
- Deterministic reduction: non-deterministic actions are replaced by deterministic actions a_1, a_2, \dots



Observation: A policy π_D for the deterministic reduction P_D is the policy for the original FOND problem P if it avoids all FOND dead-ends

→ Learns a classical general policy + state constraints

→ Encodes state constraints B by conditions C , but describing bad states

Theorem: Let \mathcal{Q}_B be a class of FOND domains. If a policy π_D for \mathcal{Q}_B avoids FOND dead-ends and state constraints B describe the dead-ends, then π_D is consisting of rules R and state constraints B is a general FOND policy for \mathcal{Q} .



International Joint Conference on Artificial Intelligence (IJCAI), Aug 5–Aug 9, 2018, Stockholm, Sweden



Propositional variables

- $Goal(s, e)$ is true if the transition (s, e) is good
- $Select(f)$ is true if the feature f is selected

Formulas

1. For every alive state s , select at least one goal transition from the self-action s :

$$\bigvee_{s' \in \text{alive}(s), f \in \text{pred}(s)} Goal(s, s')$$
2. Goal states have distance 0, i.e., for every goal state e :

$$V(e, 0)$$
3. For every alive state s , choose a goal feature:

$$Encoder \cdot \perp \vee (V(s, d))$$
4. For every good transition (s, e) , there must be one outcome leading towards the goal:

$$Goal(s, e') \wedge V(s, d') \rightarrow \bigwedge_{e' \in \text{pred}(s)} V(s', d') \rightarrow d' < d$$
5. There cannot be any good transitions to dead states, i.e., for every alive state s , dead state d :

$$-Goal(s, d) \wedge Encoder \cdot \perp$$
6. For every goal state and non-goal state s' :

$$\bigvee_{f \in \text{pred}(s'), d \in \text{pred}(s')} Select(f)$$
7. For every alive state s and dead state d' :

$$\bigvee_{f \in \text{pred}(d'), d \in \text{pred}(s)} Select(f)$$
8. Every pair of a good transition (s, d) and non-good transition (s', d') must be distinguishable by a feature:

$$Goal(s, s') \wedge -Goal(s, d) \rightarrow Distinguish(s, s', d, d')$$



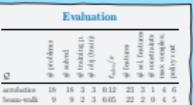
The learned policy $\pi_{enc} = \pi_{D, G}$ consists of the following rule R :

$r_1:$	$\{(d \leq 0, \neg U) \rightarrow \{d\}\}$
$r_2:$	$\{(\neg U, \neg C) \rightarrow \{U\} \mid \{d\}\}$

It has a single constraint $B = \{b\}$:

$b_1:$	$\{(B, C)\}$
--------	--------------

Proposition: The general policy π_{enc} solves the class \mathcal{Q}_{enc} of solvable acrobatics problems.



Implementation

- SAT problem solved with clingo (Gebser et al. 2011)
- Optimizations:
 - Reusing $V(s, d)$ implied by labeling safe states ($safe_i$: all children are safe, goal states always safe)
 - Only distinguish critical states from alive states (critical: dead-end but alive parent)

Come visit our poster!



Paper

Appendix

Bibliography (I)

-  Bonet, Blai, Guillem Francès, and Hector Geffner (2019). "Learning features and abstract actions for computing generalized plans". In: *AAAI*, pp. 2703–2710.
-  Bonet, Blai and Hector Geffner (2018). "Features, Projections, and Representation Change for Generalized Planning". In: *IJCAI*, pp. 4667–4673.
-  Francès, Guillem, Blai Bonet, and Hector Geffner (2021). "Learning general planning policies from small examples without supervision". In: *AAAI*, pp. 11801–11808.

Till Hofmann, Hector Geffner

IJCAI, Jeju Island, South Korea, August 6, 2024