

# Learning General Policies and Helpful Action Classifiers from Partial State Spaces

---



Dominik Drexler<sup>1</sup>



Javier Segovia-Aguas<sup>2</sup>



Jendrik Seipp<sup>1</sup>

July 23, 2022

<sup>1</sup>Linköping University, Linköping, Sweden,

<sup>2</sup>Universitat Pompeu Fabra, Barcelona, Spain

- A **general policy** encodes domain general strategies  
→ Only exist for tractable domains
- How to learn **partial policies** for **intractable domains**?
- How to learn general policies for **difficult but tractable domains**?
- Idea: find policy that solves **partial state spaces**

# Example General Policy for Visittal

## Features $\phi$

- $d$ : distance to nearest unvisited cell
- $n$ : number of unvisited cells

## Policy rules $\pi_\phi$

## Meaning

- |  |   |
|--|---|
| $\{d > 0, n > 0\} \mapsto \{d\downarrow\}$     | ; decrease distance to nearest unvisited cell |
| $\{d > 0, n > 0\} \mapsto \{d?, n\downarrow\}$ | ; visit unvisited cell                        |

- **Syntax:**

- **Policy**  $\pi_\Phi$  consists of **policy rules** of form  $C \mapsto E$  over features  $\Phi$
- For **Boolean feature**  $p$  and **numerical feature**  $n$ , we can have
  - $p, \neg p, n > 0, n = 0$  in  $C$
  - $p, \neg p, p?, n\uparrow, n\downarrow, n?$  in  $E$

- **Semantics:**

- **Transition**  $(s, a, s')$  is  $\pi_\Phi$ -**compatible** iff
  1.  $s \models C$ , and
  2.  $(s, s') \models E$
- A policy is **general** if for every alive state there exists  $\pi_\Phi$ -compatible transition
- A policy is **partial** if there exists an alive state with no  $\pi_\Phi$ -compatible transition
- In both cases, **acyclicity** over the  $\pi_\Phi$ -compatible transitions is required

# Learning Policies from Complete State Spaces

The learning problem [Francès et al., 2021]

**Given:** fully expanded state spaces  $S(P_1), \dots, S(P_n)$ , and feature pool  $\mathcal{F}$

**Find:** policy  $\pi_\Phi$  with  $\Phi \subseteq \mathcal{F}$ , s.t.,  $\pi_\Phi$  solves all  $S(P_1), \dots, S(P_n)$

- Generality of policy across whole domain can often be proven
- No provable generalization capabilities in intractable domains
- W-Max-SAT encoding: scalability problems, even more in intractable domains

→ **Idea:** observe only fragments of the state space

# Learning Policies from Partial State Spaces

## The learning problem

**Given:** partially expanded state spaces  $\mathcal{S}(P_1), \dots, \mathcal{S}(P_n)$ , and feature pool  $\mathcal{F}$

**Find:** policy  $\pi_\Phi$  with  $\Phi \subseteq \mathcal{F}$ , s.t.,  $\pi_\Phi$  solves all  $\mathcal{S}(P_1), \dots, \mathcal{S}(P_n)$

- How to obtain an informative fragment of the state space?  
→ **Idea:** optimal plans contain goal directed information

## Learning Policies from Partial State Spaces: Details

- Notation: expanded states  $\mathcal{S}$ , generated states  $\mathcal{G}$
- Initially,  $\mathcal{G}_0$  only contains initial states  $s_0$ ,  $\mathcal{S} = \emptyset$ ,  $\pi_\Phi = \emptyset$
- In each iteration  $i = 1, 2, \dots$ 
  - Sample one **optimal**  $s$ -**plan**  $p$  for  $s \in \mathcal{G}_{i-1}$  where  $\pi_\Phi$  fails
  - $\mathcal{S}_i = \mathcal{S}_{i-1} \cup \{ \text{alive states on } p \}$
  - $\mathcal{G}_i = \mathcal{G}_{i-1} \cup \{ \text{alive states on } p \} \cup \{ \text{1-step successor states } s \text{ along } p \}$
  - Generate feature pool  $\mathcal{F}_i$  with respect to  $\mathcal{G}_i$
  - Solve W-Max-SAT encoding  $\Gamma(\mathcal{S}_i, \mathcal{G}_i, \mathcal{F}_i)$  to find  $\pi_\Phi$  that solves  $\mathcal{S}_i$  **suboptimally**
  - Aim for “simplest” policy: minimize  $\sum_{f \in \Phi} \text{complexity}(f)$

## Helpful Actions with Partial Policies

- Need additional **thinking** if policy is partial
- Combine heuristic search with policy: GBFS with  $h_{FF}$  and helpful actions [Hoffmann and Nebel, 2001]
- Dual-queue approach: one prioritized queue for states reached by helpful actions
- **Action**  $a$  is  $\pi_\phi$ -**helpful** in state  $s$  if transition  $(s, a, s')$  is  $\pi_\phi$ -compatible
- **Action**  $a$  is **relaxed-helpful** in state  $s$  if  $a$  is applicable in  $s$  and part of a relaxed plan from  $s$



## Experiments: Configurations

- All configurations use GBFS and  $h_{FF}$ 
  - FF: –
  - $FF^r$ : dual-queue and **relaxed-helpful** actions
  - $FF^\pi$ : dual-queue and  $\pi_\phi$ -**helpful** actions
  - $FF_\infty$ : like FF but **greedily** following the policy after each expansion step
  - $FF_\infty^r$ : like  $FF^r$  but **greedily** following the policy after each expansion step
- Search properties: suboptimal, sound and complete

# Experiments


	FF			FF'			FF <sup>π</sup>			FF <sub>∞</sub>			FF' <sub>∞</sub>		
	S	E	T	S	E	T	S	E	T	S	E	T	S	E	T
Barman <sup>(30)</sup>	6	0.05	0.08	<b>23</b>	<b>0.50</b>	<b>0.34</b>	8	0.06	0.09	7	0.08	0.11	20	0.39	0.28
Blocks <sup>(30)</sup>	26	0.85	0.48	26	0.86	0.48	26	0.86	0.47	24	0.80	0.44	26	0.86	0.43
Blocks-clear <sup>(30)</sup>	<b>30</b>	<b>0.98</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	0.99	<b>30</b>	<b>1.00</b>	0.99
Blocks-on <sup>(30)</sup>	<b>30</b>	<b>0.95</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>0.95</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>
Childsnack <sup>(30)</sup>	1	0.00	0.01	7	0.12	0.16	3	0.02	0.03	1	0.02	0.03	<b>8</b>	<b>0.20</b>	<b>0.21</b>
Delivery <sup>(30)</sup>	<b>30</b>	<b>0.99</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>
Depots <sup>(30)</sup>	5	0.15	0.20	<b>6</b>	<b>0.18</b>	<b>0.21</b>	5	0.10	0.12	<b>6</b>	0.16	0.18	5	0.16	0.17
Driverlog <sup>(30)</sup>	8	0.13	0.16	<b>18</b>	<b>0.34</b>	<b>0.29</b>	4	0.06	0.07	6	0.10	0.10	7	0.14	0.13
Ferry <sup>(30)</sup>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>
Freecell <sup>(30)</sup>	27	0.49	0.63	26	<b>0.54</b>	<b>0.65</b>	26	0.40	0.57	26	0.48	0.62	<b>27</b>	0.52	<b>0.65</b>
Gripper <sup>(30)</sup>	<b>30</b>	0.61	0.81	<b>30</b>	0.68	0.80	<b>30</b>	0.54	0.57	<b>30</b>	<b>0.87</b>	0.94	<b>30</b>	<b>0.87</b>	0.94
Miconic <sup>(30)</sup>	<b>30</b>	<b>0.90</b>	0.97	<b>30</b>	<b>0.90</b>	<b>0.96</b>	<b>30</b>	<b>0.90</b>	0.76	<b>30</b>	<b>0.90</b>	0.84	<b>30</b>	<b>0.90</b>	0.84
N-puzzle <sup>(30)</sup>	<b>30</b>	0.88	<b>1.00</b>	<b>30</b>	0.87	<b>1.00</b>	<b>30</b>	<b>0.89</b>	<b>1.00</b>	<b>30</b>	0.88	<b>1.00</b>	<b>30</b>	0.87	<b>1.00</b>
Nomystery <sup>(30)</sup>	7	0.09	0.13	<b>10</b>	<b>0.22</b>	<b>0.26</b>	3	0.05	0.02	6	0.17	0.11	6	0.16	0.11
Parking <sup>(30)</sup>	9	0.21	0.18	11	0.29	0.25	6	0.10	0.08	<b>18</b>	<b>0.51</b>	<b>0.35</b>	17	0.49	<b>0.35</b>
Pipes-nt <sup>(30)</sup>	12	0.17	0.26	<b>24</b>	0.52	0.62	13	0.24	0.32	11	0.15	0.22	<b>24</b>	<b>0.63</b>	<b>0.71</b>
Pipes-t <sup>(30)</sup>	11	0.14	0.22	<b>28</b>	<b>0.65</b>	<b>0.63</b>	13	0.17	0.24	12	0.16	0.23	25	0.61	0.61
Reward <sup>(30)</sup>	<b>30</b>	0.92	<b>1.00</b>	<b>30</b>	0.94	<b>1.00</b>	<b>30</b>	<b>0.99</b>	0.99	<b>30</b>	<b>0.99</b>	<b>1.00</b>	<b>30</b>	<b>0.99</b>	<b>1.00</b>
Satellite <sup>(30)</sup>	10	0.33	0.36	<b>14</b>	<b>0.47</b>	<b>0.44</b>	10	0.31	0.34	10	0.33	0.36	13	0.43	0.42
Sokoban <sup>(30)</sup>	19	0.24	0.45	<b>22</b>	<b>0.26</b>	<b>0.46</b>	15	0.24	0.32	14	0.24	0.32	13	0.23	0.31
Spanner <sup>(30)</sup>	<b>0</b>	<b>0.00</b>	<b>0.00</b>	<b>0</b>	<b>0.00</b>	<b>0.00</b>	<b>27</b>	<b>0.74</b>	<b>0.13</b>	<b>30</b>	<b>0.87</b>	<b>0.30</b>	<b>30</b>	<b>0.87</b>	<b>0.30</b>
Visitall <sup>(30)</sup>	<b>5</b>	<b>0.05</b>	<b>0.11</b>	<b>5</b>	<b>0.06</b>	<b>0.11</b>	<b>6</b>	<b>0.12</b>	<b>0.14</b>	<b>25</b>	<b>0.59</b>	<b>0.44</b>	<b>25</b>	<b>0.59</b>	<b>0.44</b>
Zenotravel <sup>(30)</sup>	12	0.34	0.43	<b>15</b>	<b>0.49</b>	<b>0.67</b>	9	0.21	0.26	9	0.24	0.35	<b>15</b>	0.45	0.55
<b>Sum <sup>(690)</sup></b>	<b>398</b>	<b>0.45</b>	<b>0.50</b>	475	0.56	<b>0.58</b>	<b>414</b>	<b>0.48</b>	<b>0.46</b>	<b>445</b>	<b>0.54</b>	<b>0.52</b>	<b>501</b>	<b>0.62</b>	<b>0.58</b>

# Experiments


	FF			FF'			FF''			FF <sub>∞</sub>			FF' <sub>∞</sub>		
	S	E	T	S	E	T	S	E	T	S	E	T	S	E	T
Barman <sup>(30)</sup>	6	0.05	0.08	<b>23</b>	<b>0.50</b>	<b>0.34</b>	8	0.06	0.09	7	0.08	0.11	20	0.39	0.28
Blocks <sup>(30)</sup>	26	0.85	0.48	26	0.86	0.48	26	0.86	0.47	24	0.80	0.44	26	0.86	0.43
Blocks-clear <sup>(30)</sup>	<b>30</b>	<b>0.98</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	0.99	<b>30</b>	<b>1.00</b>	0.99
Blocks-on <sup>(30)</sup>	<b>30</b>	<b>0.95</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>0.95</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>
Childsnack <sup>(30)</sup>	1	0.00	0.01	7	0.12	0.16	3	0.02	0.03	1	0.02	0.03	<b>8</b>	<b>0.20</b>	<b>0.21</b>
Delivery <sup>(30)</sup>	<b>30</b>	<b>0.99</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>
Depots <sup>(30)</sup>	5	0.15	0.20	<b>6</b>	<b>0.18</b>	<b>0.21</b>	5	0.10	0.12	<b>6</b>	0.16	0.18	5	0.16	0.17
Driverlog <sup>(30)</sup>	8	0.13	0.16	<b>18</b>	<b>0.34</b>	<b>0.29</b>	4	0.06	0.07	6	0.10	0.10	7	0.14	0.13
Ferry <sup>(30)</sup>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>	<b>30</b>	<b>1.00</b>	<b>1.00</b>
Freecell <sup>(30)</sup>	27	0.49	0.63	26	<b>0.54</b>	<b>0.65</b>	26	0.40	0.57	26	0.48	0.62	<b>27</b>	0.52	<b>0.65</b>
Gripper <sup>(30)</sup>	<b>30</b>	0.61	0.81	<b>30</b>	0.68	0.80	<b>30</b>	0.54	0.57	<b>30</b>	<b>0.87</b>	0.94	<b>30</b>	<b>0.87</b>	0.94
Miconic <sup>(30)</sup>	<b>30</b>	<b>0.90</b>	0.97	<b>30</b>	<b>0.90</b>	<b>0.96</b>	<b>30</b>	<b>0.90</b>	0.76	<b>30</b>	<b>0.90</b>	0.84	<b>30</b>	<b>0.90</b>	0.84
N-puzzle <sup>(30)</sup>	<b>30</b>	0.88	<b>1.00</b>	<b>30</b>	0.87	<b>1.00</b>	<b>30</b>	<b>0.89</b>	<b>1.00</b>	<b>30</b>	0.88	<b>1.00</b>	<b>30</b>	0.87	<b>1.00</b>
Nomystery <sup>(30)</sup>	7	0.09	0.13	<b>10</b>	<b>0.22</b>	<b>0.26</b>	3	0.05	0.02	6	0.17	0.11	6	0.16	0.11
Parking <sup>(30)</sup>	9	0.21	0.18	11	0.29	0.25	6	0.10	0.08	<b>18</b>	<b>0.51</b>	<b>0.35</b>	17	0.49	<b>0.35</b>
Pipes-nt <sup>(30)</sup>	12	0.17	0.26	<b>24</b>	0.52	0.62	13	0.24	0.32	11	0.15	0.22	<b>24</b>	<b>0.63</b>	<b>0.71</b>
Pipes-t <sup>(30)</sup>	11	0.14	0.22	<b>28</b>	<b>0.65</b>	<b>0.63</b>	13	0.17	0.24	12	0.16	0.23	25	0.61	0.61
Reward <sup>(30)</sup>	<b>30</b>	0.92	<b>1.00</b>	<b>30</b>	0.94	<b>1.00</b>	<b>30</b>	<b>0.99</b>	0.99	<b>30</b>	<b>0.99</b>	<b>1.00</b>	<b>30</b>	<b>0.99</b>	<b>1.00</b>
Satellite <sup>(30)</sup>	10	0.33	0.36	<b>14</b>	<b>0.47</b>	<b>0.44</b>	10	0.31	0.34	10	0.33	0.36	13	0.43	0.42
Sokoban <sup>(30)</sup>	19	0.24	0.45	<b>22</b>	<b>0.26</b>	<b>0.46</b>	15	0.24	0.32	14	0.24	0.32	13	0.23	0.31
Spanner <sup>(30)</sup>	<b>0</b>	<b>0.00</b>	<b>0.00</b>	<b>0</b>	<b>0.00</b>	<b>0.00</b>	<b>27</b>	<b>0.74</b>	<b>0.13</b>	<b>30</b>	<b>0.87</b>	<b>0.30</b>	<b>30</b>	<b>0.87</b>	<b>0.30</b>
Visitall <sup>(30)</sup>	<b>5</b>	<b>0.05</b>	<b>0.11</b>	<b>5</b>	<b>0.06</b>	<b>0.11</b>	<b>6</b>	<b>0.12</b>	<b>0.14</b>	<b>25</b>	<b>0.59</b>	<b>0.44</b>	<b>25</b>	<b>0.59</b>	<b>0.44</b>
Zenotravel <sup>(30)</sup>	12	0.34	0.43	<b>15</b>	<b>0.49</b>	<b>0.67</b>	9	0.21	0.26	9	0.24	0.35	<b>15</b>	0.45	0.55
<b>Sum <sup>(690)</sup></b>	<b>398</b>	<b>0.45</b>	<b>0.50</b>	<b>475</b>	<b>0.56</b>	<b>0.58</b>	<b>414</b>	<b>0.48</b>	<b>0.46</b>	<b>445</b>	<b>0.54</b>	<b>0.52</b>	<b>501</b>	<b>0.62</b>	<b>0.58</b>

## Future Work & Conclusions

- We presented a method for combining heuristic search with partial policies learned from partial state spaces
- How to obtain more informed fragments of the state space?
- How to identify tractable fragments?

 Francès, G., Bonet, B., and Geffner, H. (2021).  
**Learning general planning policies from small examples without supervision.**

In Leyton-Brown, K. and Mausam, editors, *Proceedings of the Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI 2021)*, pages 11801–11808. AAAI Press.

 Hoffmann, J. and Nebel, B. (2001).  
**The FF planning system: Fast plan generation through heuristic search.**  
*Journal of Artificial Intelligence Research*, 14:253–302.