

Learning Strips Action Models from State-Constraints

Diego Aineto and Sergio Jiménez and Eva Onaindia

Departamento de Sistemas Informáticos y Computación

Universitat Politècnica de València.

Camino de Vera s/n. 46022 Valencia, Spain

{dieaigar,serjice,onaindia}@dsic.upv.es

Abstract

This paper extends the classical planning compilation for learning STRIPS action models from examples with the aim of learning from sets of states (no intermediate action is given). The compilation accepts state-constraints that are used to reduce the space of possible action models.

1 Introduction

Besides *plan synthesis* [?], planning action models are also useful for *plan/goal recognition* [?]. At both planning tasks, an automated planner is required to reason about action models that correctly and completely capture the possible world transitions [?]. Unfortunately, building planning action models is complex, even for planning experts, and this knowledge acquisition task is a bottleneck that limits the potential of AI planning [?].

Learning STRIPS action models is a well-studied problem with sophisticated algorithms, like ARMS [?], SLAF [?] or LOCM [?] that do not require full knowledge of the intermediate states traversed by the example plans. Motivated by recent advances on the synthesis of different kinds of generative models with classical planning [?; ?; ?], this paper introduces an innovative approach for learning STRIPS action models that can be defined as a classical planning compilation. The compilation approach is appealing by itself because opens the door to the bootstrapping of planning action models but also because:

1. Is flexible to different amounts of available input knowledge. The learning examples can range from a set of plans (with their corresponding initial and final states) to just a set of initial and final states where no intermediate state or action is observed.
2. Accepts previous knowledge about the structure of the actions in the form of partially specified action models. In the extreme, the compilation can validate whether an observed plan execution is valid for a given STRIPS action model, even if this model is not fully specified.

2 Background

This section defines the planning models used on this work and the output of the learning tasks addressed in the paper.

2.1 Classical planning

We use F to denote the set of *fluents* (propositional variables) describing a state. A *literal* l is a valuation of a fluent $f \in F$, i.e. either $l = f$ or $l = \neg f$. A set of literals L represents a partial assignment of values to fluents (WLOG we assume that L does not assign conflicting values to any fluent). We use $\mathcal{L}(F)$ to denote the set of all literal sets on F , i.e. all partial assignments of values to fluents.

A *state* s is a full assignment of values to fluents, i.e. $|s| = |F|$, so the size of the state space is $2^{|F|}$. Explicitly including negative literals $\neg f$ in states simplifies subsequent definitions but often, we will abuse notation by defining a state s only in terms of the fluents that are true in s , as is common in STRIPS planning.

A *classical planning frame* is a tuple $\Phi = \langle F, A \rangle$, where F is a set of fluents and A is a set of actions. Each action $a \in A$ comprises three sets of literals:

- $\text{pre}(a) \subseteq \mathcal{L}(F)$, called *preconditions*, the literals that must hold for the action $a \in A$ to be applicable.
- $\text{eff}^+(a) \subseteq \mathcal{L}(F)$, called *positive effects*, that defines the fluents set to true by the application of the action $a \in A$.
- $\text{eff}^-(a) \subseteq \mathcal{L}(F)$, called *negative effects*, that defines the fluents set to false by the action application.

We say that an action $a \in A$ is *applicable* in a state s iff $\text{pre}(a) \subseteq s$. The result of applying a in s is the *successor state* $\theta(s, a) = \{s \setminus \text{eff}^-(a)\} \cup \text{eff}^+(a)$.

A *classical planning problem* is a tuple $P = \langle F, A, I, G \rangle$, where I is an initial state and $G \subseteq \mathcal{L}(F)$ is a goal condition. A *plan* for P is an action sequence $\pi = \langle a_1, \dots, a_n \rangle$ that induces a state sequence $\langle s_0, s_1, \dots, s_n \rangle$ such that $s_0 = I$ and, for each $1 \leq i \leq n$, a_i is applicable in s_{i-1} and generates the successor state $s_i = \theta(s_{i-1}, a_i)$. We denote with $|\pi|$ the *plan length*. A plan π *solves* P iff $G \subseteq s_n$, i.e. if the goal condition is satisfied at the last state reached after following the application of π in I .

2.2 Classical planning with conditional effects

Our approach for learning STRIPS action models is compiling this learning task into a classical planning task with conditional effects. Conditional effects allow us to compactly define actions whose effects depend on the current state. Supporting conditional effects is now a requirement of the IPC [?]

```

(:action stack
 :parameters (?v1 ?v2 - object)
 :precondition (and (holding ?v1) (clear ?v2))
 :effect (and (not (holding ?v1))
              (not (clear ?v2))
              (handempty) (clear ?v1)
              (on ?v1 ?v2)))

```

Figure 1: STRIPS operator schema coding, in PDDL, the *stack* action from the *blocksworld*.

and many classical planners cope with conditional effects without compiling them away.

An action $a \in A$ has now a set of *preconditions* $\text{pre}(a) \in \mathcal{L}(F)$ and a set of *conditional effects* $\text{cond}(a)$. Each conditional effect $C \triangleright E \in \text{cond}(a)$ is composed of two sets of literals $C \in \mathcal{L}(F)$, the *condition*, and $E \in \mathcal{L}(F)$, the *effect*.

An action $a \in A$ is *applicable* in a state s if and only if $\text{pre}(a) \subseteq s$, and the resulting set of *triggered effects* are the effects whose conditions hold in s :

$$\text{triggered}(s, a) = \bigcup_{C \triangleright E \in \text{cond}(a), C \subseteq s} E,$$

The result of applying an action a in a state s is the *successor* state $\theta(s, a) = \{s \setminus \text{eff}_c^-(s, a)\} \cup \text{eff}_c^+(s, a)$ where $\text{eff}_c^-(s, a) \subseteq \text{triggered}(s, a)$ and $\text{eff}_c^+(s, a) \subseteq \text{triggered}(s, a)$ are the triggered *negative* and *positive* effects, respectively.

2.3 STRIPS action schemes and variable name objects

This work addresses the learning of PDDL action schemes that follow the STRIPS requirement [?; ?]. Figure 1 shows the schema, coded in PDDL, for the *stack* action from a four-operator *blocksworld* [?].

To formalize the output of the learning task, we assume that fluents F are instantiated from a set of *predicates* Ψ , as in PDDL. Each predicate $p \in \Psi$ has an argument list of arity $\text{ar}(p)$. Given a set of *objects* Ω , the set of fluents F is induced by assigning objects in Ω to the arguments of predicates in Ψ , i.e. $F = \{p(\omega) : p \in \Psi, \omega \in \Omega^{\text{ar}(p)}\}$ s.t. Ω^k is the k -th Cartesian power of Ω .

Let $\Omega_v = \{v_i\}_{i=1}^{\max_{a \in A} \text{ar}(a)}$ be a new set of objects $\Omega \cap \Omega_v = \emptyset$, denoted as *variable names*, and that is bound by the maximum arity of an action in a given planning frame. For instance, in a three-block blocksworld $\Omega = \{\text{block}_1, \text{block}_2, \text{block}_3\}$ while $\Omega_v = \{v_1, v_2\}$ because the operators with the maximum arity, *stack* and *unstack*, have two parameters each.

Let us also define F_v , a new set of fluents $F \cap F_v = \emptyset$, that results from instantiating Ψ using only the objects in Ω_v and that defines the elements that can appear in an action schema. For instance, in the blocksworld, $F_v = \{\text{handempty}, \text{holding}(v_1), \text{holding}(v_2), \text{clear}(v_1), \text{clear}(v_2), \text{ontable}(v_1), \text{ontable}(v_2), \text{on}(v_1, v_1), \text{on}(v_1, v_2), \text{on}(v_2, v_1), \text{on}(v_2, v_2)\}$.

Finally, we assume that actions $a \in A$ are instantiated from STRIPS operator schemes $\xi = \langle \text{head}(\xi), \text{pre}(\xi), \text{add}(\xi), \text{del}(\xi) \rangle$ where:

- $\text{head}(\xi) = \langle \text{name}(\xi), \text{pars}(\xi) \rangle$, is the operator *header* defined by its name and corresponding *variable names*, $\text{pars}(\xi) = \{v_i\}_{i=1}^{\text{ar}(\xi)}$. For instance, the headers for a four-operator blocksworld are: *pickup*(v_1), *putdown*(v_1), *stack*(v_1, v_2) and *unstack*(v_1, v_2).
- The preconditions $\text{pre}(\xi) \subseteq F_v$, the negative effects $\text{del}(\xi) \subseteq F_v$, and the positive effects $\text{add}(\xi) \subseteq F_v$ such that, $\text{del}(\xi) \subseteq \text{pre}(\xi)$, $\text{del}(\xi) \cap \text{add}(\xi) = \emptyset$ and $\text{pre}(\xi) \cap \text{add}(\xi) = \emptyset$.

3 Learning STRIPS action models

Learning STRIPS action models from fully available input knowledge, i.e. from plans where the *pre*- and *post*-states of every action in a plan are available, is straightforward. When any intermediate state is available, STRIPS operator schemes are derived lifting the literals that change between the pre and post-state of the corresponding action executions. Preconditions are derived lifting the minimal set of literals that appears in all the pre-states of the corresponding actions.

This section formalizes more challenging learning tasks, where less input knowledge is available:

Learning from (initial, final) state pairs.

This learning task corresponds to observing an agent acting in the world but watching only the results of its plan executions. No information about the actions in the plans is given. This learning task is formalized as $\Lambda = \langle \Psi, \Sigma \rangle$:

- Ψ is the set of predicates that define the abstract state space of a given planning domain.
- $\Sigma = \{\sigma_1, \dots, \sigma_\tau\}$ is a set of (*initial*, *final*) state pairs, that we call *labels*. Each label $\sigma_t = (s_0^t, s_n^t)$, $1 \leq t \leq \tau$, comprises the *final* state s_n^t resulting from executing an unknown plan π_t starting from the *initial* state s_0^t .

Learning from state-constraints.

A solution to Λ is a set of operator schema Ξ that is compliant just with the predicates in Ψ , and the given set of initial and final states Σ . In this learning scenario, a solution must not only determine a possible STRIPS action model but also the plans π_t , $1 \leq t \leq \tau$ that explain the given labels Σ using the learned STRIPS model. A solution to Λ' is a set of STRIPS operator schema Ξ (one schema $\xi = \langle \text{head}(\xi), \text{pre}(\xi), \text{add}(\xi), \text{del}(\xi) \rangle$ for each action with a different name in the example plans Π) compliant with the predicates in Ψ , the example plans Π , and their corresponding labels Σ . Finally a solution to Λ'' is a set of STRIPS operator schema Ξ compliant as well with the provided partially specified action model Ξ_0 .

4 Learning STRIPS action models with planning

Our approach for addressing a learning task Λ , Λ' or Λ'' , is compiling it into a classical planning task with conditional effects. The intuition behind the compilation is that a solution to the resulting classical planning task is a sequence of actions that:

```

;;; Predicates in  $\Psi$ 
(handempty) (holding ?o - object)
(clear ?o - object) (ontable ?o - object)
(on ?o1 - object ?o2 - object)

;;; Plan  $\pi_1$ 
0: (unstack A B)
1: (putdown A)
2: (unstack B C)
3: (stack B A)
4: (unstack C D)
5: (stack C B)
6: (pickup D)
7: (stack D C)

;;; Label  $\sigma_1 = (s_0^1, s_n^1)$ 

```

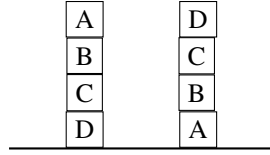


Figure 2: Example of a task for learning a STRIPS action model in the blockworld from a single labeled plan.

1. Programs the STRIPS action model Ξ . A solution plan has a *prefix* that, for each $\xi \in \Xi$, determines the fluents from F_v that belong to its $pre(\xi)$, $del(\xi)$ and $add(\xi)$ sets.
2. Validates the programmed STRIPS action model Ξ in the given input knowledge (the labels Σ , and Π and/or Ξ_0 if available). For every label $\sigma_t \in \Sigma$, a solution plan has a postfix that produces a final state s_n^t starting from the corresponding initial state s_0^t using the programmed action model Ξ . We call this process the validation of the programmed STRIPS action model Ξ , at the learning example $1 \leq t \leq \tau$.

To formalize our compilation we first define $1 \leq t \leq \tau$ classical planning instances $P_t = \langle F, \emptyset, I_t, G_t \rangle$ that belong to the same planning frame (i.e. same fluents and actions but differ in the initial state and goals). Fluents F are built instantiating the predicates in Ψ with the objects appearing in the input labels Σ . Formally $\Omega = \{o | o \in \bigcup_{1 \leq t \leq \tau} obj(s_0^t)\}$, where obj is a function that returns the set of objects that appear in a fully specified state. The set of actions, $A = \emptyset$, is empty because the action model is initially unknown. Finally, the initial state I_t is given by the state $s_0^t \in \sigma_t$ while goals G_t , are defined by the state $s_n^t \in \sigma_t$.

Now we are ready to formalize the compilations. We start with Λ , because it requires less input knowledge. Given a learning task $\Lambda = \langle \Psi, \Sigma \rangle$ the compilation outputs a classical planning task $P_\Lambda = \langle F_\Lambda, A_\Lambda, I_\Lambda, G_\Lambda \rangle$:

- F_Λ extends F with:
 - Fluents representing the programmed action model $pre_f(\xi)$, $del_f(\xi)$ and $add_f(\xi)$, for every $f \in F_v$ and $\xi \in \Xi$. If a fluent $pre_f(\xi)/del_f(\xi)/add_f(\xi)$ holds, it means that f is a precondition/negative effect/positive effect in the STRIPS operator schema $\xi \in \Xi$. For instance, the preconditions of the *stack* schema (Figure 1) are represented by fluents $pre_holding_stack.v_1$ and $pre_clear_stack.v_2$.
 - A fluent $mode_{prog}$ indicating whether the operator schemes are being programmed or validated (al-

ready programmed) and fluents $\{test_t\}_{1 \leq t \leq \tau}$, indicating the example where the action model is being validated.

- I_Λ contains the fluents from F that encode s_0^1 (the initial state of the first label), every $pre_f(\xi) \in F_\Lambda$ and $mode_{prog}$ set to true. Our compilation assumes that initially any operator schema is programmed with every possible precondition, no negative effect and no positive effect.
- $G_\Lambda = \bigcup_{1 \leq t \leq \tau} \{test_t\}$, indicates that the programmed action model is validated in all the learning examples.
- A_Λ contains actions of three kinds:
 1. Actions for *programming* an operator schema $\xi \in \Xi$:
 - Actions for **removing** a *precondition* $f \in F_v$ from the action schema $\xi \in \Xi$.

$$\begin{aligned}
pre(programPre_{f,\xi}) &= \{-del_f(\xi), \neg add_f(\xi), \\
&\quad mode_{prog}, pre_f(\xi)\}, \\
cond(programPre_{f,\xi}) &= \{\emptyset\} \triangleright \{\neg pre_f(\xi)\}.
\end{aligned}$$

- Actions for **adding** a *negative* or *positive* effect $f \in F_v$ to the action schema $\xi \in \Xi$.

$$\begin{aligned}
pre(programEff_{f,\xi}) &= \{-del_f(\xi), \neg add_f(\xi), \\
&\quad mode_{prog}\}, \\
cond(programEff_{f,\xi}) &= \{pre_f(\xi)\} \triangleright \{del_f(\xi)\}, \\
&\quad \{\neg pre_f(\xi)\} \triangleright \{add_f(\xi)\}.
\end{aligned}$$

2. Actions for *applying* an already programmed operator schema $\xi \in \Xi$ bound with the objects $\omega \subseteq \Omega^{ar(\xi)}$. We assume operators headers are known so the binding of the operator schema is done implicitly by order of appearance of the action parameters, i.e. variables $pars(\xi)$ are bound to the objects in ω appearing at the same position.

$$\begin{aligned}
pre(apply_{\xi,\omega}) &= \{pre_f(\xi) \implies p(\omega)\}_{\forall p \in \Psi, f=p(pars(\xi))}, \\
cond(apply_{\xi,\omega}) &= \{del_f(\xi)\} \triangleright \{\neg p(\omega)\}_{\forall p \in \Psi, f=p(pars(\xi))}, \\
&\quad \{add_f(\xi)\} \triangleright \{p(\omega)\}_{\forall p \in \Psi, f=p(pars(\xi))}, \\
&\quad \{mode_{prog}\} \triangleright \{\neg mode_{prog}\}.
\end{aligned}$$

3. Actions for *validating* the learning example $1 \leq t \leq \tau$.

$$\begin{aligned}
pre(validate_t) &= G_t \cup \{test_j\}_{j \in 1 \leq j < t} \\
&\quad \cup \{\neg test_j\}_{j \in t \leq j \leq \tau} \cup \{\neg mode_{prog}\}, \\
cond(validate_t) &= \{\emptyset\} \triangleright \{test_t\}.
\end{aligned}$$

Lemma 1. Any classical plan π that solves P_Λ induces an action model Ξ that solves the learning task Λ .

Proof sketch. The compilation forces that once the preconditions of an operator schema $\xi \in \Xi$ are programmed, they cannot be altered. The same happens with the positive and negative effects that define an operator schema $\xi \in \Xi$ (besides they can only be

programmed after preconditions are programmed). Once operator schemes are programmed they can only be applied because of the $mode_{prog}$ fluent. To solve P_Λ , goals $\{test_t\}$, $1 \leq t \leq \tau$ can only be achieved: executing an applicable sequence of programmed operator schemes that reaches the final state s_n^t , defined in σ_t , starting from s_0^t . If this is achieved for all the input examples $1 \leq t \leq \tau$, it means that the programmed action model Ξ is compliant with the provided input knowledge and hence, it is a solution to Λ . \square

The compilation is *complete* in the sense that it does not discard any possible STRIPS action model.

5 Constraining the learning hypothesis space with additional input knowledge

Here we show that further input knowledge can be used to constrain the space of possible action models and make the learning of STRIPS action models more practicable.

6 Evaluation

This section evaluates the performance of our approach for learning STRIPS action models starting from different amounts of available input knowledge.

Setup.

The domains used in the evaluation are IPC domains that satisfy the STRIPS requirement [?], taken from the PLANNING.DOMAINS repository [?]. We only use 5 learning examples for each domain and they are fixed for all the experiments so we can evaluate the impact of the input knowledge in the quality of the learned models. All experiments are run on an Intel Core i5 3.10 GHz x 4 with 4 GB of RAM.

Reproducibility.

We make fully available the compilation source code, the evaluation scripts and the used benchmarks at this anonymous repository <https://github.com/anonsub/strips-learning> so any experimental data reported in the paper is fully reproducible.

Planner.

The classical planner we use to solve the instances that result from our compilations is MADAGASCAR [?]. We use MADAGASCAR because its ability to deal with planning instances populated with dead-ends. In addition, SAT-based planners can apply the actions for programming preconditions in a single planning step (in parallel) because these actions do not interact. Actions for programming action effects can also be applied in a single planning step reducing significantly the planning horizon.

Metrics.

The quality of the learned models is quantified with the *precision* and *recall* metrics. Intuitively, precision gives a notion of *soundness* while recall gives a notion of the *completeness* of the learned models. Formally, $Precision = \frac{tp}{tp+fp}$, where tp is the number of true positives (predicates that correctly appear in the action model) and fp is the number of false positives (predicates appear in the learned action model that should not appear). Recall is formally defined as $Recall = \frac{tp}{tp+fn}$ where fn is the number of false negatives

(predicates that should appear in the learned action model but are missing).

7 Conclusions

References

- [Amir and Chang, 2008] Eyal Amir and Allen Chang. Learning partially observable deterministic action models. *Journal of Artificial Intelligence Research*, 33:349–402, 2008.
- [Bonet *et al.*, 2009] Blai Bonet, Héctor Palacios, and Héctor Geffner. Automatic derivation of memoryless policies and finite-state controllers using classical planners. In *ICAPS*, 2009.
- [Cresswell *et al.*, 2013] Stephen N Cresswell, Thomas L McCluskey, and Margaret M West. Acquiring planning domain models using LOCM. *The Knowledge Engineering Review*, 28(02):195–213, 2013.
- [Fox and Long, 2003] Maria Fox and Derek Long. PDDL2.1: An extension to PDDL for expressing temporal planning domains. *J. Artif. Intell. Res.(JAIR)*, 20:61–124, 2003.
- [Geffner and Bonet, 2013] Hector Geffner and Blai Bonet. A concise introduction to models and methods for automated planning, 2013.
- [Ghallab *et al.*, 2004] Malik Ghallab, Dana Nau, and Paolo Traverso. *Automated Planning: theory and practice*. Elsevier, 2004.
- [Kambhampati, 2007] Subbarao Kambhampati. Model-lite planning for the web age masses: The challenges of planning with incomplete and evolving domain models. In *Proceedings of the National Conference on Artificial Intelligence*, 2007.
- [McDermott *et al.*, 1998] Drew McDermott, Malik Ghallab, Adele Howe, Craig Knoblock, Ashwin Ram, Manuela Veloso, Daniel Weld, and David Wilkins. PDDL – The Planning Domain Definition Language, 1998.
- [Muise, 2016] Christian Muise. Planning. domains. *ICAPS system demonstration*, 2016.
- [Ramírez, 2012] Miquel Ramírez. *Plan recognition as planning*. PhD thesis, Universitat Pompeu Fabra, 2012.
- [Rintanen, 2014] Jussi Rintanen. Madagascar: Scalable planning with sat. *Proceedings of the 8th International Planning Competition (IPC-2014)*, 2014.
- [Segovia-Aguas *et al.*, 2016] Javier Segovia-Aguas, Sergio Jiménez, and Anders Jonsson. Hierarchical finite state controllers for generalized planning. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, pages 3235–3241. AAAI Press, 2016.
- [Segovia-Aguas *et al.*, 2017] Javier Segovia-Aguas, Sergio Jiménez, and Anders Jonsson. Generating context-free grammars using classical planning. In *International Joint Conference on Artificial Intelligence*, 2017.
- [Slaney and Thiébaux, 2001] John Slaney and Sylvie Thiébaux. Blocks world revisited. *Artificial Intelligence*, 125(1-2):119–153, 2001.
- [Vallati *et al.*, 2015] Mauro Vallati, Lukáš Chrpá, Marek Grzes, Thomas L McCluskey, Mark Roberts, and Scott Sanner. The 2014 international planning competition: Progress and trends. *AI Magazine*, 36(3):90–98, 2015.
- [Yang *et al.*, 2007] Qiang Yang, Kangheng Wu, and Yunfei Jiang. Learning action models from plan examples using weighted max-sat. *Artificial Intelligence*, 171(2-3):107–143, 2007.