

Learning STRIPS action models with classical planning

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 presents a novel approach for learning STRIPS action models from examples that compiles this particular inductive learning task into classical planning. Interestingly, the compilation approach 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 (no intermediate action or state is given). What is more, the compilation accepts partially specified action models and can be used to validate whether certain observations of plan executions follow a given STRIPS action model, even if this model is not fully specified.

Introduction

Besides *plan synthesis* (Ghallab, Nau, and Traverso 2004; Geffner and Bonet 2013), planning action models are also useful for *plan/goal recognition* (Ramírez 2012). 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 automated planning (Kambhampati 2007).

On the other hand, Machine Learning (ML) has shown to be able to compute a wide range of different kinds of models from examples (Michalski, Carbonell, and Mitchell 2013). The application of inductive ML to the learning of STRIPS action models (Fikes and Nilsson 1971), the vanilla action model for planning, is not straightforward though:

- The inputs to ML algorithms (the learning/training data) usually are finite vectors encoding the value of fixed features in a given set of objects. The input for learning planning action models are observations of plan executions (where each plan possibly has different length).
- The traditional output of ML algorithms is a scalar value (an integer, in the case of *classification* tasks, or a real value, in the case of *regression* tasks). When learning STRIPS action models the output is, for each action, the

sets of preconditions, negative and positive effects, that define the possible state transitions.

The learning of STRIPS action models is a well-studied problem with sophisticated algorithms, like ARMS (Yang, Wu, and Jiang 2007), SLAF (Amir and Chang 2008) or LOCM (Cresswell, McCluskey, and West 2013) 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 (Bonet, Palacios, and Geffner 2009; Segovia-Aguas, Jiménez, and Jonsson 2016; 2017), 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 because it 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 could be used to validate whether observed plan executions are valid for a given STRIPS action model.

The paper is organized as follows. The first section presents the classical planning model, its extension to conditional effects (which is a requirement of the proposed compilation) and formalizes the STRIPS action model (the output of the addressed learning task). The second section formalizes the learning of STRIPS action models with regard to different amounts of available input knowledge. The third and fourth sections describe our approach for addressing this particular learning task by compiling it into classical planning and how previous knowledge can be introduced in the compilation to constrain the space of possible action models and make learning more practicable. Finally, last sections report the data collected during the empirical evaluation of our approach, discuss the strengths and weaknesses of our approach and propose several opportunities for future research.

Background

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

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 total 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 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:

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

We say that an action $a \in A$ is *applicable* in state s iff $\text{pre}(a) \subseteq s$, and 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 $G \in \mathcal{L}(F)$ 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 $s_i = \theta(s_{i-1}, a_i)$. 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 .

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. Many classical planners cope with conditional effects without compiling them away. In fact, supporting conditional effects is a requirement of IPC-2014 (Vallati et al. 2015) and IPC-2018.

Now an action $a \in A$ has a set of literals $\text{pre}(a) \in \mathcal{L}(F)$, called the *precondition*, 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 state s if and only if $\text{pre}(a) \subseteq s$, and the resulting set of *triggered effects* is

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

```
(:action stack
:parameters (?x1 ?x2 - object)
:precondition (and (holding ?x1) (clear ?x2))
:effect (and (not (holding ?x1)) (not (clear ?x2))
(handempty) (clear ?x1) (on ?x1 ?x2)))
```

Figure 1: Example of a STRIPS operator schema that corresponds to the *stack* action from the *blocksworld* represented in PDDL.

i.e. effects whose conditions hold in s .

The result of applying action a in a state s is the *successor state* $\theta(s, a) = \{s \setminus \text{del}(s, a) \cup \text{add}(s, a)\}$ where $\text{del}(s, a) \subseteq \text{triggered}(s, a)$ and $\text{add}(s, a) \subseteq \text{triggered}(s, a)$ are the triggered negative and positive effects.

The STRIPS action schema and the variable objects

This work addresses the learning of PDDL action schemes that follow the STRIPS requirement (McDermott et al. 1998; Fox and Long 2003). Figure 1 shows the schema that corresponds to the *stack* action from a four-operator *blocksworld* (Slaney and Thiébaux 2001).

To formalize the output of the learning task, we assume that there exists a set of *predicates* Ψ and that fluents F are instantiated from these predicates, 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 then 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 . Likewise, we assume that actions $a \in A$ are instantiated from STRIPS operator schemes.

Let $\Omega_v = \{v_i\}_{i=1}^{\max_{a \in A} \text{ar}(a)}$ be a new set of objects, $\Omega \cap \Omega_v = \emptyset$, that represents *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 define also a new set of fluents F_v , s.t. $F \cap F_v = \emptyset$, that results instantiating Ψ using only *variable objects* in Ω_v . This set defines the elements that can appear in an action schema. In the blocksworld $F_v = \{\text{handempty}, \text{holding}(v_1), \text{holding}(v_2), \text{clear}(v_1), \text{clear}(v_2), \text{ontable}(v_1), \text{ontable}(v_2), \text{on}(v_1, v_1), \text{on}(v_1, v_2), \text{on}(v_2, v_1), \text{on}(v_2, v_2)\}$.

We are now ready to define Ξ , a set of operator schema $\xi = \langle \text{head}(\xi), \text{pre}(\xi), \text{add}(\xi), \text{del}(\xi) \rangle$ such that:

- $\text{head}(\xi) = \langle \text{name}(\xi), \text{pars}(\xi) \rangle$, is the operator *header* defined by its name and $\text{pars}(\xi) = \{v_i\}_{i=1}^{\text{ar}(\xi)}$, an enumeration of the variable objects bound by the operator arity. The headers for a 4-operator blocksworld are: $\text{pickup}(v_1)$, $\text{putdown}(v_1)$, $\text{stack}(v_1, v_2)$ and $\text{unstack}(v_1, v_2)$.
- 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 $\text{pre}(\xi) \cap \text{add}(\xi) = \emptyset$.

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. In this case, because intermediate states are 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 that correspond to the same operator schema.

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 however, we assume that the headers of the operators schema are known.

This learning task is formalized as $\Lambda = \langle \Psi, \Sigma \rangle$:

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

A solution to the learning task Λ is a set of operator schema Ξ that is compliant 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 that explain the given states Σ using the learned STRIPS model.

Learning from labeled plans

Here we augment the amount of provided input knowledge and define the learning task as $\Lambda' = \langle \Psi, \Sigma, \Pi \rangle$.

Now $\Pi = \{\pi_1, \dots, \pi_\tau\}$ is a given set of example plans where each plan $\pi_t = \langle a_1^t, \dots, a_n^t \rangle$, $1 \leq t \leq \tau$, is an action sequence that induces the corresponding state sequence $\langle s_0^t, s_1^t, \dots, s_n^t \rangle$ such that, for each $1 \leq i \leq n$, a_i^t is applicable in the state s_{i-1}^t and generates the successor state $s_i^t = \theta(s_{i-1}^t, a_i^t)$.

A solution to a learning task Λ' is a set of STRIPS operator schema Ξ (with 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) that is compliant with the predicates in Ψ , the example plans Π , and their corresponding labels Σ .

Figure 2 shows an example of a learning task Λ' in the blockworld (the figure shows the content of Ψ , Π and Σ). This learning task has a single learning example, $\Pi = \{\pi_1\}$ and $\Sigma = \{\sigma_1\}$, that corresponds to observing the execution of a plan for inverting a four-block tower.

;;; Predicates in Ψ

```
(handempty) (holding ?o - object)
(clear ?o - object) (ontable ?o - object)
(on ?o1 - object ?o2 - object)
```

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

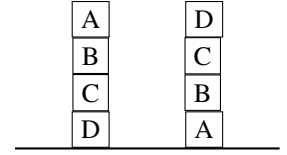


Figure 2: Example of a task for learning a STRIPS action model in the blockworld.

Learning STRIPS action models with planning

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

1. Programs the STRIPS action model Ξ . For each $\xi \in \Xi$, the solution plan has a prefix that determines the literals that belong to its *pre*(ξ), *del*(ξ) and *add*(ξ) sets.
2. Validates the programmed STRIPS action model Ξ using the given input knowledge (the set of labels Σ , and Π if available). For every label $\sigma_t \in \Sigma$, then Ξ is used to produce a final state s_n^t starting from the corresponding initial state s_0^t . We call this process the validation of the programmed STRIPS action model Ξ , at the learning example $1 \leq t \leq \tau$. If information about the corresponding plan $\pi_t \in \Pi$ is available, then it is also used in the validation.

To formalize our compilations 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 and differ only in the initial state and goals). The fluents F are built instantiating the predicates in Ψ with the objects appearing in the inputs. Formally $\Omega = \{o|o \in \bigcup_{1 \leq t \leq \tau} \text{obj}(s_0^t)\}$, 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 this learning task requires fewer input knowledge. Given a learning task $\Lambda = \langle \Psi, \Sigma \rangle$ the compilation outputs a classical planning task $P_\Lambda = \langle F_\Lambda, A_\Lambda, I_\Lambda, G_\Lambda \rangle$ such that:

- F_Λ extends F with:
 - Fluents representing the programmed action model $\text{pre}_f(\xi)$, $\text{del}_f(\xi)$ and $\text{add}_f(\xi)$, for every $f \in F_v$ and $\xi \in \Xi$. If a fluent $\text{pre}_f(\xi)/\text{del}_f(\xi)/\text{add}_f(\xi)$ holds, it means that f is a precondition/negative effect/positive effect in the STRIPS operator schema $\xi \in \Xi$.

For instance in the working example the preconditions of the *stack* schema are represented by the fluents *pre_holding_stack_v1* and *pre_clear_stack_v2* set to true.

- A fluent *mode_{prog}* indicating whether the operator schemes are being programmed or they are being validated (already programmed).
- Fluents $\{test_t\}_{1 \leq t \leq \tau}$, indicating the learning example where the programmed action model is being validated.
- I_Λ , contains the fluents from F that encode s_0^1 (the initial state of the first learning example), every $pre_f(\xi) \in F_\Lambda$ (our compilation assumes that initially any operator schema is programmed with every possible precondition, no negative effect and no positive effect) and *mode_{prog}* set to true.
- $G_\Lambda = \bigcup_{1 \leq t \leq \tau} \{test_t\}$, indicates that the programmed action model is validated in all the learning examples.
- A_Λ contains actions of three different kinds:
 1. The actions for programming an operator schema $\xi \in \Xi$, which includes:
 - The actions for **removing** a *precondition* $f \in F_v$ from the action schema $\xi \in \Xi$.

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

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

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

2. The actions for applying an already programmed operator schema $\xi \in \Xi$ bound with the objects $\omega \subseteq \Omega^{ar(\xi)}$. The binding of the operator schema is done implicitly by order of appearance of the action parameters, i.e. the variables in *vars*(ξ) are bound to the objects in ω appearing at the same position. Figure 3 shows the PDDL encoding of the action for applying a programmed operator *stack*.

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

3. The actions for updating the learning example $1 \leq t \leq \tau$ where the programmed action model is validated.

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

```
(:action apply_stack
:parameters (?o1 - object ?o2 - object)
:precondition
  (and (or (not (pre_on_stack_v1_v1)) (on ?o1 ?o1))
        (or (not (pre_on_stack_v1_v2)) (on ?o1 ?o2))
        (or (not (pre_on_stack_v2_v1)) (on ?o2 ?o1))
        (or (not (pre_on_stack_v2_v2)) (on ?o2 ?o2))
        (or (not (pre_ontable_stack_v1)) (ontable ?o1))
        (or (not (pre_ontable_stack_v2)) (ontable ?o2))
        (or (not (pre_clear_stack_v1)) (clear ?o1))
        (or (not (pre_clear_stack_v2)) (clear ?o2))
        (or (not (pre_holding_stack_v1)) (holding ?o1))
        (or (not (pre_holding_stack_v2)) (holding ?o2))
        (or (not (pre_handempty_stack)) (handempty)))
:effect
  (and (when (del_on_stack_v1_v1) (not (on ?o1 ?o1)))
        (when (del_on_stack_v1_v2) (not (on ?o1 ?o2)))
        (when (del_on_stack_v2_v1) (not (on ?o2 ?o1)))
        (when (del_on_stack_v2_v2) (not (on ?o2 ?o2)))
        (when (del_ontable_stack_v1) (not (ontable ?o1)))
        (when (del_ontable_stack_v2) (not (ontable ?o2)))
        (when (del_clear_stack_v1) (not (clear ?o1)))
        (when (del_clear_stack_v2) (not (clear ?o2)))
        (when (del_holding_stack_v1) (not (holding ?o1)))
        (when (del_holding_stack_v2) (not (holding ?o2)))
        (when (del_handempty_stack) (not (handempty)))
        (when (add_on_stack_v1_v1) (on ?o1 ?o1))
        (when (add_on_stack_v1_v2) (on ?o1 ?o2))
        (when (add_on_stack_v2_v1) (on ?o2 ?o1))
        (when (add_on_stack_v2_v2) (on ?o2 ?o2))
        (when (add_ontable_stack_v1) (ontable ?o1))
        (when (add_ontable_stack_v2) (ontable ?o2))
        (when (add_clear_stack_v1) (clear ?o1))
        (when (add_clear_stack_v2) (clear ?o2))
        (when (add_holding_stack_v1) (holding ?o1))
        (when (add_holding_stack_v2) (holding ?o2))
        (when (add_handempty_stack) (handempty))
        (when (modeProg) (not (modeProg)))))
```

Figure 3: Action for applying an already programmed operator schema *stack* as encoded in PDDL.

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 the preconditions are programmed). Furthermore, once operator schemes are programmed they can only be applied because of the *mode_{prog}* fluent. To solve P_Λ , there is only one way of achieving goals $\{test_t\}$, $1 \leq t \leq \tau$: 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 (defined in this same label). If this is achieved for all the input examples $1 \leq t \leq \tau$ of the learning task, 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.

```

0 : (program_pre_clear_stack_v1)
1 : (program_pre_handempty_stack)
2 : (program_pre_holding_stack_v2)
3 : (program_pre_on_stack_v1_v1)
4 : (program_pre_on_stack_v1_v2)
5 : (program_pre_on_stack_v2_v1)
6 : (program_pre_on_stack_v2_v2)
7 : (program_pre_ontable_stack_v1)
8 : (program_pre_ontable_stack_v2)
9 : (program_eff_clear_stack_v1)
10 : (program_eff_clear_stack_v2)
11 : (program_eff_handempty_stack)
12 : (program_eff_holding_stack_v1)
13 : (program_eff_on_stack_v1_v2)
14 : (apply_unstack a b i1 i2)
15 : (apply_putdown a i2 i3)
16 : (apply_unstack b c i3 i4)
17 : (apply_stack b a i4 i5)
18 : (apply_unstack c d i5 i6)
19 : (apply_stack c b i6 i7)
20 : (apply_pickup d i7 i8)
21 : (apply_stack d c i8 i9)
22 : (validate_1)

```

Figure 4: Example of a plan for programming and validating the operator schema *stack* using the plan π_1 and label σ_1 shown in Figure 2 as well as previously specified operator schemes for *pickup*, *putdown* and *unstack*.

Figure 4 shows a classical plan for solving a learning task of the Λ kind. This classical plan programs and validates the operator schema *stack* from the blocksworld (previously specified operator schemes for *pickup*, *putdown* and *unstack* are given), using the plan π_1 and label σ_1 shown in Figure 2. The plan steps [0, 8] are the actions for programming the preconditions of the *stack* operator, steps [9, 13] are the actions for programming the operator effects and steps [14, 22] are actions for validating the programmed operators. The details of the compilation are explained in the remaining of the section.

Constraining the learning hypothesis space with additional input knowledge

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

Example plans

Here we extend our compilation to address the learning scenario where, a set of plans is available. Given a learning task $\Lambda = \langle \Psi, \Pi, \Sigma \rangle$, the compilation outputs a classical planning task $P_\Lambda = \langle F_\Lambda, A_\Lambda, I_\Lambda, G_\Lambda \rangle$ that extends P'_Λ as follows:

- F_Λ extends $F_{\Lambda'}$ with $F_\Pi = \{plan(name(\xi), \Omega^{ar(\xi)}, j)\}$, the fluents to code the j steps of plans in Π , where $F_{\pi_t} \subseteq F_\Pi$ encodes $\pi_t \in \Pi$. Fluents at_j and $next_{j,j_2}$, $1 \leq j < j_2 \leq n$, are also added to represent the current and next plan steps.
- I_Λ is extended with the fluents from F_Π that encode the plan $\pi_1 \in \Pi$, the fluents at_1 and $\{next_{j,j_2}\}$, $1 \leq j < j_2 \leq$

n , for indicating where to start validating the programmed action model. Goals are the same as in the previous compilation $G_\Lambda = G'_\Lambda = \bigcup_{1 \leq t \leq \tau} \{test_t\}$.

- With respect to A_Λ .

1. The actions for programming the preconditions/effects of a given operator schema $\xi \in \Xi$ are the same.
2. The actions for applying an already programmed operator have an extra precondition $f \in F_\Pi$, that encodes the current plan step, and extra conditional effects $\{at_j\} \triangleright \{\neg at_j, at_{j+1}\}_{\forall j \in [1, n]}$ for advancing to the next plan step. This mechanism forces that these actions are only applied as in the example plans.
3. The actions for updating the active test have an extra precondition, $at_{|\pi_t|}$, to indicate that the current plan π_t was fully executed and extra conditional effects to load the next plan π_{t+1} :

$$\{f\} \triangleright \{\neg f\}_{f \in F_{\pi_t}}, \{\emptyset\} \triangleright \{f\}_{f \in F_{\pi_{t+1}}}, \{\emptyset\} \triangleright \{\neg at_{|\pi_t|}, at_1\}.$$

Previously specified action models

Our compilation does not require to start learning STRIPSaction models from scratch and accepts partially specified operator schema where some preconditions and effects are a priori known.

The known preconditions and effects are encoded as fluents $pre_f(\xi)$, $del_f(\xi)$ and $add_f(\xi)$ that are part of the initial state $I_{\Lambda'}$. The corresponding programming actions ($programPre_{f,\xi}$ and $programEff_{f,\xi}$) become unnecessary and are removed from $A_{\Lambda'}$ making the classical planning task $P_{\Lambda'}$ easier to be solved. To illustrate this, the classical plan of Figure 4 is a solution to a learning task for getting the blocksworld action model where the schemes for *pickup*, *putdown* and *unstack* are known in advance.

In the extreme, when a fully specified STRIPSaction model is given, the compilation validates whether an observed plan follows the given model. In this case, if a solution plan is found to $P_{\Lambda'}$, it means that the given STRIPSaction model is *valid* for the given set of examples. If $P_{\Lambda'}$ is unsolvable it means that the given STRIPSaction model is invalid since it is not compliant with all the given examples. Tools for plan validation like VAL (Howey, Long, and Fox 2004) could also be used at this point.

Static predicates

A *static predicate* $p \in \Psi$ is a predicate that does not appear in the effects of any action schema (Fox and Long 1998). Therefore, one can get rid of the mechanism for programming static predicates as action effects keeping the compilation complete. Formally, given a static predicate p :

- Fluents $del_f(\xi)$ and $add_f(\xi)$, such that $f \in F_v$ is an instantiation of the static predicate p in the set of *variable* objects Ω_v , can be discarded for every $\xi \in \Xi$.
- Actions $programEff_{f,\xi}$ (s.t. $f \in F_v$ is an instantiation of p in Ω_v) can also be discarded for every $\xi \in \Xi$.

Static predicates can also be used to constrain the space of possible preconditions by looking at the given labels. If a precondition $f \in F_v$ (s.t. $f \in F_v$ is an instantiation

of a static predicate in Ω_v) is not possible according to Σ , then the fluents $pre_f(\xi)$ and actions $programPre_{f,\xi}$ can be discarded for every $\xi \in \Xi$. For instance in the *zenotravel* domain fluents $pre_next_board.v1.v1$, $pre_next_debark.v1.v1$, $pre_next_fly.v1.v1$, $pre_next_zoom.v1.v1$, $pre_next_refuel.v1.v1$ can be discarded (as well as their corresponding programming actions) because a precondition ($next \text{ ?v1 ?v1 - flevel}$) cannot be compliant with any state in Σ .

Likewise, static predicates can constrain the space of possible preconditions looking at the given example plans. In this case, fluents $pre_f(\xi)$ and actions $programPre_{f,\xi}$ are discarded for every $\xi \in \Xi$ if a precondition $f \in F_v$ (s.t. $f \in F_v$ is an instantiation of a static predicate in Ω_v) is not possible according to Π . Back to the *zenotravel* domain, if a example plan $\pi_t \in \Pi$ contains the action ($fly \text{ plane1 city2 city0 f13 f12}$) and the corresponding label $\sigma_t \in \Sigma$ contains the static fluent ($next \text{ f12 f13}$) but it does not include ($next \text{ f12 f12}$), ($next \text{ f13 f13}$) or ($next \text{ f13 f12}$) it means that the only possible precondition including the static predicate is $pre_next_fly.v1.v2$.

In the evaluation of our approach we do not assume that the set of static predicates is given but compute a set of candidates of static predicates from the input examples to the learning task.

Evaluation

We evaluated our approach for learning STRIPS action models from different amounts of available input knowledge.

Experimental setup

All the experiments run on an Intel Core i5 3.10 GHz x 4 with 4 GB of RAM. The domains used in the evaluation are IPC domains that satisfy the STRIPS requirement (Fox and Long 2003), and are taken from the PLANNING.DOMAINS repository (Muise 2016).

The classical planner. The classical planner used to solve the instances that result from our compilations is the SAT-based classical planner MADAGASCAR (Rintanen 2014). We decided to use a SAT-based planner because its ability to deal with planning instances populated with dead-ends. In addition, in a SAT-based classical planner all the actions for programming the STRIPS actions preconditions can be done in parallel, in a single planning step, because they do not interact reducing significantly the planning horizon. Likewise, all the actions for programming an action effect can also be applied in a single planning step.

The evaluation metrics. The quality of the learned models is quantified with the *precision* and *recall* metrics. These two metrics are frequently used in pattern recognition, information retrieval and binary classification and are more informative than simply counting the number of model errors or computing the *symmetric difference* between the learned model and the baseline (Davis and Goadrich 2006).

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). On the other hand 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).

The learning examples.

Learning action models from example plans

Here we assess the performance of our learning approach when addressing Λ learning tasks. Here the *precision* and *recall* metrics are computed in the learned models with respect to the actual models taken as reference. Table 1 summarizes the obtained results, precision (**P**) and recall (**R**) is computed separately for the preconditions (**Pre**), positive effects (**Add**) and negative Effects (**Del**).

	Pre		Add		Del		All	
	P	R	P	R	P	R	P	R
Blocks	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Driverlog	1.0	0.43	0.67	0.86	1.0	0.86	0.89	0.71
Ferry	0.80	0.57	1.0	1.0	1.0	1.0	0.93	0.86
Floortile	0.54	0.64	0.90	0.82	1.0	0.91	0.81	0.79
Gripper	1.0	0.67	1.0	1.0	1.0	1.0	1.0	0.89
Miconic	0.75	0.33	1.0	0.75	0.75	1.0	0.83	0.69
Satellite	0.67	0.29	1.0	1.0	1.0	0.75	0.89	0.68
Transport	1.0	0.30	0.57	0.80	1.0	0.60	0.86	0.57
Visitall	1.0	0.50	1.0	1.0	1.0	1.0	1.0	0.83
Zenotravel	1.0	0.36	0.75	0.86	1.0	0.71	0.92	0.64
Total	0.88	0.51	0.89	0.91	0.98	0.89	0.91	0.77

Table 1: The *precision* and *recall* values achieved by the learned models in the Λ learning tasks.

Table 2 reports the total planning time, and the preprocessing time (in seconds) invested by the classical planner MADAGASCAR to solve the classical planning tasks that result from our compilation as well as the plan length of the found solutions.

Domain	Total time	Preprocess	Plan length
Blocks	0.04	0.00	72
Driverlog	0.10	0.06	127
Ferry	0.20	0.10	58
Floortile	2.15	1.21	170
Gripper	0.01	0.00	45
Miconic	0.10	0.09	54
Satellite	0.20	0.14	67
Transport	0.20	0.18	60
Visitall	0.43	0.39	37
Zenotravel	0.37	0.35	70

Table 2: Planning time, preprocessing time and Plan length.

Exploiting static predicates

Here we repeat the previous evaluation but compute now a candidates of static predicates from the inputs to the learning task and exploit this knowledge to reduce the space of possible action models. The set of considered candidate *static*

predicates is the set of predicates whose instantiation appear the same in the initial and final states for all the learning examples $1 \leq t \leq \tau$. Tables 3 and 4 show the obtained results.

	Pre		Add		Del		All	
	P	R	P	R	P	R	P	R
Blocks	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Driverlog	1.0	0.5	0.75	0.86	1.0	0.71	0.91	0.69
Ferry	0.80	0.57	1.0	1.0	1.0	1.0	0.93	0.86
Floortile	0.71	0.78	0.90	0.82	1.0	0.91	0.87	0.83
Gripper	1.0	0.67	1.0	1.0	1.0	1.0	1.0	0.89
Miconic	0.88	0.78	1.0	0.75	0.75	1.0	0.88	0.84
Satellite	0.85	0.79	0.83	1.0	1.0	0.75	0.89	0.85
Transport	1.0	0.60	0.80	0.80	1.0	0.60	0.93	0.67
Visitall	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
Zenotravel	1.0	0.64	0.75	0.86	1.0	0.71	0.92	0.74
Total	0.92	0.73	0.90	0.91	0.98	0.87	0.93	0.84

Table 3: The *precision* and *recall* values achieved by the learned models in the Λ learning tasks exploiting the computed static predicates.

Domain	Total time	Preprocess	Plan length
Blocks	0.03	0.00	72
Driverlog	0.12	0.07	102
Ferry	0.25	0.15	58
Floortile	0.73	0.43	80
Gripper	0.01	0.00	45
Miconic	0.06	0.04	40
Satellite	0.18	0.12	60
Transport	0.18	0.18	44
Visitall	0.43	0.42	33
Zenotravel	0.26	0.25	53

Table 4: Planning time, preprocessing time and Plan length exploiting the computed static predicates.

Learning from partially specified action models

Here we evaluate the ability of the compilation to support partially specified action models. In this case instead of learning the action models from scratch the model of half of the actions is given as input of the learning task. Tables 5 and 6 summarizes the obtained results.

	Pre		Add		Del	
	P	R	P	R	P	R
Blocks	1.0	1.0	1.0	1.0	1.0	1.0
Driverlog	1.0	0.4	1.0	1.0	1.0	0.7
Ferry	1.0	0.6	0.5	1.0	1.0	1.0
Floortile	0.5	0.6	1.0	0.8	1.0	0.8
Gripper	1.0	0.5	1.0	1.0	1.0	1.0
Miconic	1.0	0.2	1.0	1.0	1.0	1.0
Satellite	1.0	0.5	0.5	1.0	1.0	1.0
Transport	1.0	0.2	1.0	1.0	1.0	0.5
Zenotravel	1.0	0.5	0.6	0.6	1.0	1.0

Table 5: The *precision* and *recall* values achieved by the learned models in the different domains.

Domain	Total time	Preprocess	Plan length
Blocks	0.07	0.01	54
Driverlog	0.06	0.03	71
Ferry	0.18	0.14	50
Floortile	0.81	0.68	92
Gripper	0.01	0.00	37
Miconic	0.06	0.04	40
Satellite	0.18	0.12	43
Transport	0.15	0.14	33
Zenotravel	2.54	2.54	42

Table 6: Planning time, preprocessing time and Plan length.

Learning action models from example states

When the set of input planning task is not available, the planner determines also the actions that must satisfying the input labels so actions can reformulated and still be compliant with the learning examples. For instance a blocksworld can be learned where the operator `stack` is defined with the preconditions and effects of the `unstack` operator and vice versa.

Past and future of learning action models

Back in the 90's various systems were aimed at learning operators mostly via interaction with the environment. The proposal of the LIVE system was to capture and formulate observable features of objects and use them to acquire and refine operators (Shen and Simon 1989). OBSERVER, on the other hand, updated preconditions and effects by removing and adding facts, respectively, accordingly to observations (Wang 1995). These early works were based on direct lifting of the observed states and supported by exploratory plans or external teachers, but none provided a theoretical justification for this second source of knowledge. The recent work reported in (Walsh and Littman 2008) shows that although learning STRIPS operators from pure interaction with the environment requires an exponential number of samples, access to an external teacher can provide solution traces on demand.

Whilst the aforementioned works deal with full state observability, action model learning has also been studied in domains where there is partial or missing state observability (Jiménez et al. 2012). The ARMS approach (Yang, Wu, and Jiang 2007), which works when no or partial intermediate states are given, extracts frequent action sets from plans that share a common set of parameters and frequent literal-action pairs with the help of the initial and goal state of the plans. The extracted information is used then to define a set of weighted constraints that must hold for all the plans to be correct, and solves the weighted propositional satisfiability problem with a MAX-SAT solver. Since complex planning problems may result in a too large MAX-SAT representation to be solved efficiently, ARMS implements a hill-climbing method that models the actions approximately. Consequently, the output model of ARMS may be inconsistent with the examples. The approach presented in (Amir and Chang 2008) deals with partial observability. Given a formula representing the initial belief state, a sequence of

executed actions and the corresponding partially observed states, it builds a complete explanation of observations by models of actions through a CNF formula. The learning algorithm updates the formula of the belief state with every action and observation in the sequence. This update makes sure that the new formula represents all the transition relations consistent with the actions and observations. The formula returned at the end includes all consistent models, which can then be retrieved with additional processing. Unlike the previous approaches, the one described in (Mourão et al. 2012) deals with both missing and noisy predicates in the observations. An action model is first learnt by constructing a set of kernel classifiers which tolerate noise and partial observability and then STRIPSrules are derived from the classifiers' parameters.

The approach for learning planning models that most likely works with the least information possible is LOCM (Cresswell, McCluskey, and West 2013), where the only required input are the example training plans without the need for providing the initial and goal state information. LOCM exploits assumptions about the kind of domain model it has to generate. Particularly, it assumes a domain consists of a collection of objects (called *sorts*) where each object has a defined set of states that can be captured by a parameterized Finite State Machine (FSM). The induction of FSM draws upon assumptions like the continuity of object transitions or the association of parameters between consecutive actions in the training sequence. The intuitive assumptions of LOCM yield a learning model heavily reliant on the kind of domain structure and not capable of properly deriving domain theories in which the state of a sort is subject to different state machines. This limitation is later overcome by LOCM2 by forming separate state machines for a sort, each containing a subset of the full transition set for the sort (Cresswell and Gregory 2011). The last contribution of the LOCM family, LOP (LOCM with Optimized Plans (Gregory and Cresswell 2016)), addresses the problem of inducing static predicates in action models by identifying the set of minimal static relations associated to an operator. Because LOCM models induce similar models for domains with similar structure (e.g. Blocksworld and Freecell), they face problems at identifying models that do not have static relations (Blocksworld) and those that do have static predicates (Freecell). In order to mitigate this drawback, LOP applies a post-processing step after the LOCM analysis which requires additional information about the plans, namely a set of optimal plans to be used in the learning phase.

=====

(Stern and Juba 2017).

Conclusions

The paper presented a novel approach for learning STRIPSaction models from examples exclusively using classical planning. The approach is flexible to different amounts of available input knowledge and accepts partially specified action models.

The empirical evaluation shows that, when example plans are available, our approach can compute accurate action models from small sets of learning examples and investing

very small learning times (less than a second time in most of the domains). When action plans are not available our approach is still able to produce action models compliant with the input information. In this case, since learning is not constrained by actions it can change the semantics of the operators. An interesting research direction related to this issue is *domain reformulation* to use actions in a more efficient way, reduce the set of actions identifying dispensable information or exploiting features that allow more compact solutions like the *reachable* or *movable* features in the *Sokoban* domain (Ivankovic and Haslum 2015).

The size of the compiled classical planning instances depends on the number of input examples. The empirical results show that our approach can generate non-trivial models from very small data sets. In addition the SAT-based planner MADAGASCAR is particularly suitable for the approach because its ability to deal with planning instances populated with dead-ends and because all the actions for programming the STRIPSactions preconditions can be done in parallel, in a single planning step, because they do not interact reducing significantly the planning horizon. The same happens with all the actions for programming an action effect can also be applied in a single planning step.

As far as we know, this is the first work on learning action models using a classical planner. This novel approach opens the door to the bootstrapping of action model that is, a planner collecting by its own the necessary data for learning an action model for a given planning task. Generating *informative* examples for learning planning action models is still a challenging open issue. Planning actions include preconditions that are only satisfied by specific sequences of actions, and often, with a low probability of being chosen by chance (Fern, Yoon, and Givan 2004). The success of recent algorithms for exploring planning tasks (Francés et al. 2017) motivates the development of novel techniques that autonomously collect the learning examples.

Acknowledgment

Diego Aineto is partially supported by the *FPU* program funded by the Spanish government. Sergio Jiménez is partially supported by the *Ramon y Cajal* program funded by the Spanish government.

References

- Amir, E., and Chang, A. 2008. Learning partially observable deterministic action models. *Journal of Artificial Intelligence Research* 33:349–402.
- Bonet, B.; Palacios, H.; and Geffner, H. 2009. Automatic derivation of memoryless policies and finite-state controllers using classical planners. In *ICAPS*.
- Cresswell, S., and Gregory, P. 2011. Generalised domain model acquisition from action traces. In *ICAPS*.
- Cresswell, S. N.; McCluskey, T. L.; and West, M. M. 2013. Acquiring planning domain models using locm. *The Knowledge Engineering Review* 28(02):195–213.
- Davis, J., and Goadrich, M. 2006. The relationship between precision-recall and roc curves. In *Proceedings of the 23rd international conference on Machine learning*, 233–240. ACM.

- Fern, A.; Yoon, S. W.; and Givan, R. 2004. Learning domain-specific control knowledge from random walks. In *ICAPS*, 191–199.
- Fikes, R. E., and Nilsson, N. J. 1971. Strips: A new approach to the application of theorem proving to problem solving. *Artificial intelligence* 2(3-4):189–208.
- Fox, M., and Long, D. 1998. The automatic inference of state invariants in tim. *Journal of Artificial Intelligence Research* 9:367–421.
- Fox, M., and Long, D. 2003. Pddl2. 1: An extension to pddl for expressing temporal planning domains. *J. Artif. Intell. Res.(JAIR)* 20:61–124.
- Francés, G.; Ramrez, M.; Lipovetzky, N.; and Geffner, H. 2017. Purely declarative action representations are overrated: Classical planning with simulators. In *International Joint Conference on Artificial Intelligence*.
- Geffner, H., and Bonet, B. 2013. A concise introduction to models and methods for automated planning.
- Ghallab, M.; Nau, D.; and Traverso, P. 2004. *Automated Planning: theory and practice*. Elsevier.
- Gregory, P., and Cresswell, S. 2016. Domain model acquisition in the presence of static relations in the LOP system. In *International Joint Conference on Artificial Intelligence, IJCAI-16*, 4160–4164.
- Howey, R.; Long, D.; and Fox, M. 2004. Val: Automatic plan validation, continuous effects and mixed initiative planning using pddl. In *Tools with Artificial Intelligence, 2004. ICTAI 2004. 16th IEEE International Conference on*, 294–301. IEEE.
- Ivankovic, F., and Haslum, P. 2015. Optimal planning with axioms. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.
- Jiménez, S.; De La Rosa, T.; Fernández, S.; Fernández, F.; and Borrajo, D. 2012. A review of machine learning for automated planning. *The Knowledge Engineering Review* 27(04):433–467.
- Kambhampati, S. 2007. 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*.
- McDermott, D.; Ghallab, M.; Howe, A.; Knoblock, C.; Ram, A.; Veloso, M.; Weld, D.; and Wilkins, D. 1998. Pddl-the planning domain definition language.
- Michalski, R. S.; Carbonell, J. G.; and Mitchell, T. M. 2013. *Machine learning: An artificial intelligence approach*. Springer Science & Business Media.
- Mourão, K.; Zettlemoyer, L. S.; Petrick, R. P. A.; and Steedman, M. 2012. Learning STRIPS operators from noisy and incomplete observations. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 614–623.
- Muise, C. 2016. Planning domains. *ICAPS system demonstration*.
- Ramírez, M. 2012. *Plan recognition as planning*. Ph.D. Dissertation, Universitat Pompeu Fabra.
- Rintanen, J. 2014. Madagascar: Scalable planning with sat. *Proceedings of the 8th International Planning Competition (IPC-2014)*.
- Segovia-Aguas, J.; Jiménez, S.; and Jonsson, A. 2016. Hierarchical finite state controllers for generalized planning. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, 3235–3241. AAAI Press.
- Segovia-Aguas, J.; Jiménez, S.; and Jonsson, A. 2017. Generating context-free grammars using classical planning. In *International Joint Conference on Artificial Intelligence*.
- Shen, W., and Simon, H. A. 1989. Rule creation and rule learning through environmental exploration. In *International Joint Conference on Artificial Intelligence*, 675–680.
- Slaney, J., and Thiébaux, S. 2001. Blocks world revisited. *Artificial Intelligence* 125(1-2):119–153.
- Stern, R., and Juba, B. 2017. Efficient, safe, and probably approximately complete learning of action models. *IJCAI*.
- Vallati, M.; Chrapa, L.; Grzes, M.; McCluskey, T. L.; Roberts, M.; and Sanner, S. 2015. The 2014 international planning competition: Progress and trends. *AI Magazine* 36(3):90–98.
- Walsh, T. J., and Littman, M. L. 2008. Efficient learning of action schemas and web-service descriptions. In *Proceedings of the National Conference on Artificial Intelligence*, 714–719.
- Wang, X. 1995. Learning by observation and practice: An incremental approach for planning operator acquisition. In *In Proceedings of the 12th International Conference on Machine Learning*, 549–557.
- Yang, Q.; Wu, K.; and Jiang, Y. 2007. Learning action models from plan examples using weighted max-sat. *Artificial Intelligence* 171(2-3):107–143.