

Incremental Learning of Discrete Planning Domains from Continuous Perceptions

Luciano Serafini, Paolo Traverso

Fondazione Bruno Kessler

{serafini, traverso}@fbk.eu

paper adhering to the theme on **model acquisition**

Abstract

We propose a framework for learning discrete deterministic planning domains. In this framework, an agent learns the domain by observing the action effects through continuous features that describe the state of the environment after the execution of each action. Besides, the agent learns its *perception function*, i.e., a probabilistic mapping between state variables and sensor data represented as a vector of continuous random variables called *perception variables*. We define an algorithm that updates the planning domain and the perception function by (i) introducing new states, either by extending the possible values of state variables, or by weakening their constraints; (ii) adapts the perception function to fit the observed data (iii) adapts the transition function on the basis of the executed actions and the effects observed via the perception function. The framework is able to deal with exogenous events that happen in the environment.

1 Introduction and Motivations

Automated Planning methods and techniques rely on models of the world, usually called Planning Domains. The (automated) acquisition of these models is widely recognised as a challenging bottleneck, see, e.g., the KEPS workshops and the ICKEPS competition.¹ The automated learning of planning domains is a way to address this challenge. Indeed, most often, it is impossible to specify a complete and correct model of the world. Moreover, most of the times a model needs to be updated and adapted to a changing environment.

Several and different learning approaches have been proposed so far. Some works on domain model acquisition focus on the problem of learning action schema, see, e.g. [Gregory and Cresswell, 2016; McCluskey *et al.*, 2009; Cresswell *et al.*, 2013; Mourão *et al.*, 2012; Mehta *et al.*, 2011; Zhuo and Yang, 2014]. Learning planning operators and domain models from plan examples and solution traces [Yang *et al.*, 2007; Zhuo *et al.*, 2010; Zhuo and Kambhampati, 2013;

¹The Knowledge Engineering for Planning and Scheduling (KEPS) Workshop and Competition (ICKEPS)

Henaff *et al.*, 2017] and learning probabilistic planning operators have also been investigated [Pasula *et al.*, 2004; Zettlemoyer *et al.*, 2005; Pasula *et al.*, 2007].

We propose a framework in which a discrete deterministic planning domain is extended with a perception function, i.e., a probabilistic mapping between state variables and observations from the real world represented by continuous variables, called perception variables. The perception function is represented by a conditional probability distribution that computes the likelihood of observing some values of the perception variables given an assignment to state variables.

We define an algorithm that builds an abstract deterministic finite planning domain and a perception function by executing actions and observing the effects through perception variables. The only information about the real world that is available to the learning algorithm is provided by the perceptions variables. The algorithm does not have access to a continuous model of the dynamics of the world. In several cases, such model is not available or is too difficult to provide.

The learning algorithm can start either “from scratch” (i.e., with an “empty planning domain”), or from some prior knowledge expressed with an initial discrete planning domain and perception function. The algorithm incrementally learns the values of the state variables, the description of the transition function, the constraints on state variables, and the perception function. The framework provides the ability to learn and adapt to unexpected situations, i.e., some constraints on state variables have been violated, or the domain of some state variables should be extended with new values.

The paper is structured as follows. Section 2 formalises the planning domain, including the perception function. Section 3 defines the incremental learning algorithm. In Section 4 we show how the algorithm works with an explanatory example that shows the potentialities of the framework. We finally discuss related work, conclusions, and future work.

2 Perceived Planning Domains

A (*deterministic*) *planning domain* is a triple $\mathcal{D} = \langle S, A, \gamma \rangle$, composed of a finite non empty set of states S , a finite non empty set of actions A , and a state transition function $\gamma : S \times A \rightarrow S$. Each state $s \in S$ is represented with a vector of *state variables* ranging over a finite set of values. Let $\mathbf{V} = \langle V_1, \dots, V_m \rangle$ be a vector of m state variables. Let $\mathbf{D} = \{D_1, \dots, D_k\}$ be a set of non empty finite sets,

called *domains*. Let \mathbf{Dom} be a function that assigns a domain $\mathbf{Dom}(V)$ to each variable V of \mathbf{V} . The set $\mathbf{Dom}(V)$ is the set of values that can be assigned to the variable V . For every $\mathbf{W} \subseteq \mathbf{V}$, we use $\mathbf{Dom}(\mathbf{W})$ to denote the cross product of the domains of all the variables in \mathbf{W} , namely $\mathbf{Dom}(\mathbf{W}) = \prod_{V \in \mathbf{W}} \mathbf{Dom}(V)$. For every $\mathbf{w} \in \mathbf{Dom}(\mathbf{W})$, we use $\mathbf{W} = \mathbf{w}$ to denote the (partial) assignment to each variable $V \in \mathbf{W}$ to $v \in \mathbf{Dom}(V)$. If \mathbf{W} is the entire set of variables \mathbf{V} then $\mathbf{V} = \mathbf{v}$ is a *total assignment*. A state $s \in S$ is a total assignment, i.e., a set of assignments that assigns a value $v \in \mathbf{Dom}(V)$ to every state variable V . We use $s[V]$ to denote the value assigned by s to V . Not every total assignment necessarily corresponds to a state. The set of states S of a planning domain is a subset of the total assignments. S can be specified with a set of *constraints* between values of state variables. For instance, the fact that V and V' must take different values can be represented by the constraint $V \neq V'$. In this paper we suppose that constraints are expressed using propositional combination (via \wedge , \vee and \neg) of the atomic proposition $V = v$, and $V = V'$, for $V, V' \in \mathbf{V}$ and $v \in \mathbf{Dom}(V)$.

We assume that the transition function γ is specified with action language, resulting in a compact representation. In this paper we adopt a simple action language, which specifies γ through a set of rules of the form

$$r : \text{prec}(r) \xrightarrow{a} \text{eff}(r) \quad (1)$$

where $a \in A$, $\text{prec}(r)$ is a propositional formula in the language of the constraints, and $\text{eff}(r)$ is a partial assignment $\mathbf{V} = \mathbf{v}'$. For every action a and state s , $s' = \gamma(a, s)$ is the state obtained after the execution of a in s , and is defined as

$$s'[V] = \begin{cases} v & \text{if } \exists r \text{ for } a, \text{ such that } s \models \text{prec}(r) \text{ and} \\ & \text{eff}(r) \text{ contains } V = v \\ s[V] & \text{otherwise} \end{cases}$$

In order to guarantee that $\gamma(a, s)$ is deterministic, we impose that for every pair of rules r and r' , defining the action a , we have that if $\text{prec}(r) \wedge \text{prec}(r')$ is consistent then $\text{eff}(r) \cup \text{eff}(r')$ does not contain $V = v_1$ and $V = v_2$ for $v_1 \neq v_2$.

The agent perceives the world through a vector $\mathbf{X} = \langle X_1, \dots, X_n \rangle$ of continuous variables ranging over real numbers, called *perception variables*. A *perception function*, is a function $f : \mathbb{R}^n \times \mathbf{Dom}(\mathbf{V}) \rightarrow \mathbb{R}^+$, such that for every $\mathbf{x} \in \mathbb{R}^n$ and total assignment $\mathbf{V} = \mathbf{v}$, $f(\mathbf{x}, \mathbf{v}) = p(\mathbf{x} \mid \mathbf{V} = \mathbf{v})$, where $p(\mathbf{x} \mid \mathbf{V} = \mathbf{v})$ is a probability density function (PDF) that can be factorised as follows:

$$p(\mathbf{x} \mid \mathbf{V} = \mathbf{v}) = \prod_{i=1}^n p_{X_i}(x_i \mid \mathbf{V}_{J_i} = \mathbf{v}_{j_i})$$

where \mathbf{V}_{J_i} is a subset of the state variables \mathbf{V} .

Definition 1 (Extended planning domain) An extended planning domain is a pair $\langle \mathcal{D}, f \rangle$ where \mathcal{D} is a planning domain and f a perception function on the states of \mathcal{D} .

Hereafter, if not explicitly specified, with “planning domain” we will refer to extended planning domain.

Example 1 (The “Robot-Pack-Cat (RPC) Flat”) The RPC-Flat is composed of 6 rooms (named from A to F), see Figure 1. In this flat there are a robot, a pack, and a cat. The robot can move from one room to adjacent rooms, load, transport and unload the pack. The cat moves around randomly and can also jump on top of the robot. The robot is equipped with an RFID reader able to perceive the presence in the room of the pack, which is equipped with a proximity sensor tag. Suppose that the robot has only partial knowl-

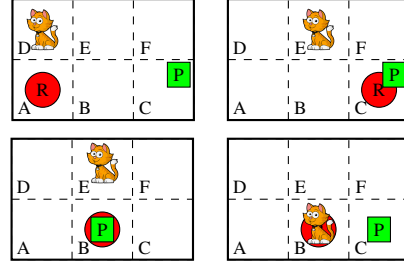


Figure 1: Four possible situations in the RPC-flat

edge about the flat and its dynamics. It believes that there are only 4 rooms (ignoring the room C and F), it ignores also the presence of the cat. The robot represents its partial knowledge with the following planning domain: The states are represented by three state variables: $\text{loc}(r)$, $\text{loc}(p)$, and loaded , which represent the position of the robot, the position of the pack, and whether the robot is loaded. There are two domains i.e., $\mathbf{D} = \{\text{room}, \text{nr_of_carried_objects}\}$ where $\text{room} = \{0, 1, 2, 3\}$ and $\text{nr_of_carried_objects} = \{0, 1\}$, with $\mathbf{Dom}(\text{loc}(r)) = \text{room}$, $\mathbf{Dom}(\text{loc}(p)) = \text{room}$, and $\mathbf{Dom}(\text{loaded}) = \text{nr_of_carried_objects}$. Notice that the robot assumes that there are only 4 rooms and 1 object to be carried.

Not all the state variable assignments are states (in S), indeed, when the robot is carrying the pack, their position must be the same. This can be formalized by the constraint:

$$\text{loaded} = 1 \rightarrow \text{loc}(r) = \text{loc}(p) \quad (2)$$

The set A of actions include N, S, E, W (that stand for the robot moves north, south, east, and west, respectively), L, and U (that stand for the robot loads and unloads the pack). Examples of a specification for E and L are the following:

$$\begin{aligned} \text{loc}(r) = 0 &\xrightarrow{E} \text{loc}(r) = 1 \\ \text{loc}(r) = 0 \wedge \text{loaded} = 1 &\xrightarrow{E} \text{loc}(p) = 1 \\ \text{loc}(r) = \text{loc}(p) &\xrightarrow{L} \text{loaded} = 1 \end{aligned}$$

The robot has the following perception variables:

- X, Y with $\mathbf{Dom}(X) = \mathbf{Dom}(Y) = \mathbb{R}$ are the x - and y -coordinates of the position of the robot;
- T with $\mathbf{Dom}(T) = [0, 1]$ is the output of RFID reader. If the pack and the robot are in the same room then the value of T is close to 1, otherwise it is close to 0;
- W with $\mathbf{Dom}(W) = \mathbb{R}^+$ is the weight currently carried by the robot.

The perception function is factorized as follows:

$$p(x, y, z, w \mid \text{loc}(r), \text{loc}(p), \text{loaded}) = p_X(x \mid \text{loc}(r)) \cdot p_Y(y \mid \text{loc}(r)) \cdot p_T(t \mid \text{loc}(r), \text{loc}(p)) \cdot p_W(w \mid \text{loaded}) \text{ where:}$$

$$\begin{aligned} p_X(x \mid \text{loc}(r)) &= \mathcal{N}(x \mid \mu_{X, \text{loc}(r)}, \sigma) \\ \mu_{X, \text{loc}(r)} &= \text{loc}(r) \bmod 2 + 0.5, \\ p_Y(y \mid \text{loc}(r)) &= \mathcal{N}(y \mid \mu_{Y, \text{loc}(r)}, \sigma) \\ \mu_{Y, \text{loc}(r)} &= \text{loc}(r) \div 2 + 0.5, \\ p_T(t \mid \text{loc}(r), \text{loc}(p)) &= \text{B}(t \mid \alpha_{\text{loc}(r), \text{loc}(p)}, \beta_{\text{loc}(r), \text{loc}(p)}) \\ \alpha_{\text{loc}(r), \text{loc}(p)} &= \cdot \mathbb{1}_{\text{loc}(r)=\text{loc}(p)} + 2 \cdot \mathbb{1}_{\text{loc}(r) \neq \text{loc}(p)} \\ \beta_{\text{loc}(r), \text{loc}(p)} &= 2 \cdot \mathbb{1}_{\text{loc}(r)=\text{loc}(p)} + 1 \cdot \mathbb{1}_{\text{loc}(r) \neq \text{loc}(p)} \\ p_W(w \mid \text{loaded}) &= \Gamma(w \mid k_{\text{loaded}}, \theta_{\text{loaded}}) \\ k_{\text{loaded}} &= \text{loaded} + 1, \theta = 1 \end{aligned}$$

3 The Incremental Learning Algorithm

The *Acting and Learning Planning-domains* algorithm, ALP, described in Algorithm 1, not only learns/updates the transitions of a planning domain, but it can also learn/update the perception function, and extend the set of states, either by weakening some constraints, or by extending the domains of some state variables. ALP can start “from scratch”, i.e., from the simplest planning domain, where each variable domain $D \in \mathbf{D}$ is equal to $\{0\}$, without constraints, and an empty γ . Alternatively, ALP can start from any non empty planning domain corresponding to some “prior knowledge” about the world.

Given a planning domain with a set of state variables in input, ALP requires the perception function $f(\mathbf{x}, \mathbf{v}) = \prod_{i=1}^n p_{X_i}(\cdot \mid \mathbf{V}_{J_i} = \mathbf{v}_{J_i})$ to be defined for all variable assignments $\mathbf{V} = \mathbf{v}$. Furthermore, since ALP introduces new values in the domain of state variables when the perception function of a perceived value \mathbf{x} is too low, we need a method to initialise the perception function for these new values. For this reason ALP requires in input also an initialiser p_{init, X_i} , for every perception variable X_i , that returns a PDF for any observation \mathbf{x} . Moreover, ALP requires in input some additional *update parameters*, α , β , γ , and δ , all in $[0, 1]$, which determine how much the agent trusts in the various components of the model. In this section, we will explain the meaning of each parameter.

ALP iteratively refines the current planning domain \mathcal{D} with the associated perception function f , by executing the actions proposed by EXPLORE (line 4),² and by observing the action effects through the perception variables \mathbf{x} (line 5). In order to determine the next state s'_0 (from line 6 to line 15), ALP firstly computes ABOVETHRESHOLD(\mathbf{x}, S) for the observation \mathbf{x} , which corresponds to the set of states such that the likelihood of observing each x_i is above the threshold $(1 - \epsilon) \cdot \max p_{X_i}$. Formally: ABOVETHRESHOLD(\mathbf{x}, S) returns the set $\{s \in S \mid \forall i, p_{X_i}(x_i \mid s[\mathbf{V}_{J_i}]) \geq (1 - \epsilon) \cdot \max p_{X_i}\}$. Intuitively, ABOVETHRESHOLD selects a set of states that are

² A naïve implementation of EXPLORE can be a random generator of actions. A smarter strategy can take into account how much has the agent already learned, which portion of the domain has been already explored, and the part that still requires more learning.

Algorithm 1 ALP

Require: $\mathcal{D} = \langle S, A, \gamma \rangle$ {Initial planning domain}
Require: $f = \prod p_{X_i}$ {Initial perception function}
Require: s_0 {Initial state}
Require: $\alpha, \beta, \delta, \epsilon$ {Update parameters}
Require: p_{init, X_i} {Perception initialization for X_i }
Require: MAXITER {Maximum number of exploration steps}
1: $T \leftarrow \langle \rangle$ {The empty history of transitions}
2: $O \leftarrow \langle \rangle$ {The empty history of observations}
3: **for** ITER $\leftarrow 1$ **to** MAXITER **do**
4: $a \leftarrow \text{EXPLORE}(\mathcal{D}, s_0)$
5: $\mathbf{x} \leftarrow \text{ACT}(a)$
6: $S'_0 \leftarrow \text{ABOVETHRESHOLD}(\mathbf{x}, S)$
7: **if** $S'_0 = \emptyset$ **then**
8: $S'_0 \leftarrow \text{ABOVETHRESHOLD}(\mathbf{x}, \text{Dom}(\mathbf{V}) \setminus S)$
9: **if** $S'_0 = \emptyset$ **then**
10: $\mathbf{D} \leftarrow \text{EXTENDDOM}(\mathbf{D}, f, \mathbf{x})$
11: $f \leftarrow \text{EXTENDF}(\mathbf{D}, f, \mathbf{x})$
12: $S'_0 \leftarrow \text{ABOVETHRESHOLD}(\mathbf{x}, \text{Dom}_{\text{new}}(\mathbf{V}))$
13: **end if**
14: **end if**
15: $s'_0 \leftarrow \text{ONEOF}(\arg\max_{s \in S'_0} f(\mathbf{x}, s) \cdot \text{sim}(s, \gamma(s_0, a) \mid \delta))$
16: **if** $s'_0 \notin S$ **then**
17: $S \leftarrow S \cup \{s'_0\}$
18: **end if**
19: $T \leftarrow \text{APPEND}(T, \langle s_0, \pi(s_0), s'_0 \rangle)$ {extend the transition history with the last one}
20: $O \leftarrow \text{APPEND}(O, \langle s'_0, \mathbf{x} \rangle)$ {extend the observation history with the last one}
21: $\gamma \leftarrow \text{UPDATETRANS}(\gamma, T \mid \alpha)$
22: $f \leftarrow \text{UPDATEPERC}(f, O \mid \beta)$
23: $s_0 \leftarrow s'_0$
24: **end for**

the candidates to be the next state, i.e., those states for which the likelihood of observing x_i is higher than a certain threshold defined by the parameter $\epsilon \in [0, 1]$. At one extreme, when $\epsilon = 1$, ABOVETHRESHOLD(\mathbf{x}, S) selects all states in S . On the other extreme, if $\epsilon = 0$, ABOVETHRESHOLD(\mathbf{x}, S) selects only those states in which $f(\mathbf{x}, s)$ reaches its maximum value. The lower ϵ , the higher chance to introduce new states. Intuitively, ϵ expresses how much we believe that the set of states learned so far are sufficient for the planning domain to model the real world.

At line 7, if there are no assignments among the current states that pass the threshold, then ALP considers the assignments which are not in the set of states, i.e., $\text{Dom}(\mathbf{V}) \setminus S$ (line 8). If, even in this case, ABOVETHRESHOLD returns the empty set (line 9), then we need to extend the possible assignments to variable by extending their domain. This is performed by EXTENDDOM (line 10), which extends the domain of one or more state variable.

EXTENDDOM (see Algorithm 2) takes in input the set \mathbf{D} of current state variables domains, the perception function f , and the current observation \mathbf{x} . It starts by selecting one assignment s that maximises the likelihood of observing \mathbf{x} . Then it computes the set $X_{<lik}$ of perception variables X_i where the likelihood of the perceived value x_i w.r.t. the state s is below the threshold (line 2). For every variable X_i in $X_{<lik}$, EXTENDDOM selects a domain in $\mathbf{D}_{J_i} = \{\text{Dom}(V) \mid$

Algorithm 2 EXTENDDOM

Require: D {The set of domain of state variables}
Require: $f = \prod p_{X_i}$ {Perception function}
Require: \mathbf{x} {The result of a perception}
1: $s \leftarrow \text{ONEOF}(\text{argmax}_{s \in \text{Dom}(\mathbf{V})} f(\mathbf{x}, s))$
2: $\mathbf{X}_{<lik} \leftarrow \{X_i \mid p_{X_i}(x_i \mid s[\mathbf{V}_{J_i}]) < (1 - \epsilon) \cdot \max p_{X_i}\}$
3: $\mathbf{D}_H \leftarrow$ minimal hitting set of $\{\mathbf{D}_{J_i}\}_{X_i \in \mathbf{X}_{<lik}}$
4: **for** $D \in \mathbf{D}_H$ **do**
5: $D \leftarrow D \cup \{|D|\}$
6: **end for**
7: **return** D

Algorithm 3 EXTENDF

Require: D {The set of domain of state variables}
Require: \mathbf{x} {The result of a perception}
Require: p_{init, X_i} {Perception initializer for X_i }
1: **for** $\mathbf{v} \in \text{Dom}(\mathbf{V})$ **do**
2: **for** $X_i \in \mathbf{X}$ **do**
3: **if** $p_{X_i}(\cdot \mid \mathbf{V}_{J_i} = \mathbf{v}_{J_i})$ is not defined **then**
4: $p_{X_i}(\cdot \mid \mathbf{V}_{J_i} = \mathbf{v}_{J_i}) = p_{init, X_i}(x_i)$
5: **end if**
6: **end for**
7: **end for**

$V \in \mathbf{V}_{J_i}$ to be extended with a new value. Since we want to minimize the number of values introduced, we choose to extend the set of domains \mathbf{D}_H that is a minimal hitting set³ for $\{\mathbf{D}_{J_i}\}_{X_i \in \mathbf{X}_{<lik}}$. (line 3). Each domain in $D \in \mathbf{D}_H$ is extended with a new value, resulting in the set of $|D| + 1$ elements $\{0, 1, 2, \dots, |D|\}$ (line 5).

After executing EXTENDDOM, ALP calls EXTENDF (line 11) to initialise the perception function for the newly introduced states. EXTENDF (see Algorithm 3) does this for all the variables without perception function (line 4). The introduction of the new values for state variables, and the initialisation guarantees that ABOVETHRESHOLD returns a non empty set S'_0 of assignments. Then ALP selects the next state s'_0 among the elements of S'_0 (line 15). The next state is one among the states that maximize the product of the likelihood of observing \mathbf{x} and the similarity with the state predicted by the transition function learned so far, i.e., $\gamma(a, s_0)$. Ideally the next state will be the one that maximises the likelihood of the perceived values, and the closest to the state predicted by the model. These two sources of information however could be contradictory, therefore we have to jointly maximize their product.

The similarity/distance measure, $\text{sim}(s, s' \mid \delta)$ for $s, s' \in \text{Dom}(\mathbf{V})$ is defined as

$$\prod_{i=1}^m \frac{1 + \delta \cdot (\mathbb{1}_{s[V_i]=s'(V_i)} \cdot (|\text{Dom}(V_i)| - 1) - \mathbb{1}_{s[V_i] \neq s'(V_i)})}{1 + \delta(|\text{Dom}(V_i)| - 1)}$$

The parameter $\delta \in [0, 1]$ allows us to adjust the similarity measure between states. At one extreme, if $\delta = 0$, then every

³A set of A is an hitting set of a family of sets $\{B_i\}_{i=1}^n$ if $A \cap B_i \neq \emptyset$ for every i . A is a minimal hitting set if there is no hitting set A' for $\{B_i\}_{i=1}^n$ with $|A'| < |A|$.

state is similar to every other state, i.e., $\text{sim}(s, s' \mid \delta) = 1$, and the similarity does not play any role in the maximisation. If $\delta = 1$, sim coincides with the equality relation, i.e., $\text{sim}(s, s' \mid \delta) = \mathbb{1}_{s=s'}$, which implies that the maximization will always return $\gamma(a, s_0)$. The interesting case is when $\delta \in (0, 1)$. The lower δ , the more we trust in the perceptions of the agent's sensors. The higher δ , the more we trust in the model learned so far.

If s'_0 is not part of the current set of states S , we have to include it by weakening the constraints. Let C_1, \dots, C_k be the set of constraints defining S . To specify $S \cup \{s'_0\}$, we have to weaken each C_i as follows

$$C_i \vee \bigwedge_{V \in \mathbf{V}} V = s'_0[V] \quad (3)$$

and if the new values v_{new} is introduced we have to add the following constraint:

$$V = v_{new} \rightarrow \bigwedge_{V' \neq V} V' = s'_0[V'] \quad (4)$$

for every variable V for which the domain $\text{Dom}(V)$ has been extended with the new value v_{new} .

Proposition 3.1 *Let $\{C'_i\}_{i=1}^h$ be the set of constraints resulting from the revision of $\{C_i\}_{i=1}^k$ according to the rules (3) and (4), then $s \models \bigwedge_{i=1}^h C'_i$ if and only if $s \in S \cup \{s'_0\}$.*

ALP then extends the sequence of transitions T and of observations O , and learn the new transition function γ and the new perception function f . The functions UPDATETRANS and UPDATEPERC update the transition function γ and the perception function f , respectively, depending on the data available in T and O . The update functions take into account (i) the current model, (ii) what has been observed in the past, i.e., T and O , and (iii) what has been just observed, i.e., $\langle s_0, a, s'_0 \rangle$ and $\langle s'_0, \mathbf{x} \rangle$. The update functions can be defined in several different ways, depending on whether we follow a cautious strategy, where changes are made only if there is a certain number of evidences from acting and perceiving the real world, or a more impulsive reaction to what the agent has just observed. In the following, we describe in detail how we create/update transitions, and how we create/update perception functions.

Updating transitions. UPDATETRANS decides whether and how to update the transition function. If s'_0 is the state that maximises the product of the perception function and of the similarity, and s'_0 is different from the state predicted by the planning domain, i.e., $s'_0 \neq \gamma(a, s_0)$, then γ may need to be revised to take into account this discrepancy. Since our domain is deterministic (the transition γ must lead to a single state), if the execution of an action leads to an unexpected state, we have only two options: either change γ with the new transition or not. We propose the following transition update function that depends on α : We define $\text{UPDATETRANS}(\gamma, T)(s, a) = s'$ where s' is a state that maximizes

$$\alpha \cdot \mathbb{1}_{s'=\gamma(s,a)} + (1 - \alpha) \cdot |\{i \mid T_i = \langle s, a, s' \rangle\}| \quad (5)$$

where T_i is the i -th element of T , and $\alpha \in [0, 1]$. Notice that, if $\alpha = 1$, we are extremely cautious, we strongly believe in our model of the world, and we never change the transition γ . Conversely, if $\alpha = 0$, we are extremely impulsive,

we do not trust our model, and just one evidence makes us to change the model. In the intermediate cases, $\alpha \in (0, 1)$, depending on the value of α , we need more or less evidence to change the planning domain. In order to update γ , we have to revise the action specifications. We replace every rule r about a of the form $r : prec(r) \xrightarrow{a} eff(r)$, such that $s \models prec(r)$ and not $s \models eff(r)$ with the following rules for every V_i

$$r'_i : prem(r) \wedge V_i \neq s[V_i] \xrightarrow{a} eff(r)$$

and the following rule for all j , such that $s[V_j] \neq s'[V_j]$

$$r''_j : \bigwedge_{j=1}^m V_j = s[V_j] \xrightarrow{a} V_j = s'[V_j]$$

Notice that this method might generate a proliferation of very specific rules. Therefore after this step it is convenient to apply some algorithm for rule factorisation. Examples of factorization rules are the following:

$$\begin{aligned} \Gamma, V = v \xrightarrow{a} V' = v' \\ \Gamma, V \neq v \xrightarrow{a} V' = v' \end{aligned} \text{ are merged in } \Gamma \xrightarrow{a} V' = v'$$

Another example, is the following. Suppose that $Dom(V) = \{0, 1, 2\}$, then:

$$\begin{aligned} \Gamma, V = 0 \xrightarrow{a} V' = v' \\ \Gamma, V = 1 \xrightarrow{a} V' = v' \end{aligned} \text{ are merged in } \Gamma, V \neq 2 \xrightarrow{a} V' = v'$$

A final example is the following. Suppose that $Dom(V)$ and $Dom(V')$ are equal to $\{0, 1\}$ then

$$\begin{aligned} \Gamma, V = 0, V' = 0 \xrightarrow{a} V'' = v'' \\ \Gamma, V = 1, V' = 1 \xrightarrow{a} V'' = v'' \end{aligned} \text{ are merged in } \Gamma, V = V' \xrightarrow{a} V'' = v''$$

Dealing with rule factorisation can be considered as a separate topic and for lack of space is not treated in this paper. However, this operation results crucial in order to generate compact and “understandable” description of the transition function.

Updating the perception function. The update of the perception function is based on the current perception function $f(\mathbf{x}, s)$ for $s \in S$ and the set of observations O . We suppose that each function p_{X_i} composing the perception function $f = \prod_i p_{X_i}$, belongs to a parametric family with parameters θ_{X_i} . For every partial assignment $\mathbf{V}_{J_i} = \mathbf{v}_{J_i}$ to the state variables \mathbf{V}_{J_i} from which X_i depends on, $p_{X_i}(\cdot | \mathbf{v}_{J_i})$ is obtained by setting the parameters θ_{X_i} to some value $\theta_{X_i, \mathbf{v}_{J_i}}$. In Example 1, p_X is a Gaussian distribution with parameters $\theta_X = \langle \mu_X, \sigma_X \rangle$. For every value $r \in Dom(\text{loc}(r))$, $\mu_{X,r}$ and $\sigma_{X,r}$ are the mean and the standard deviation of p_X and $p_X(x | \mu_{X,r}, \sigma_{X,r}) = \mathcal{N}(x, \mu = \mu_{X,r}, \sigma = \sigma_{X,r})$ expresses the likelihood of observing x when the robot is in the room r . We denote by $\theta_{\mathbf{X}}$ all the parameters in $\theta_{X_1}, \dots, \theta_{X_n}$, and $\theta_{\mathbf{X}, \mathbf{v}}$, for $\mathbf{v} \in Dom(\mathbf{V})$, their instantiations $\theta_{X_1, \mathbf{v}_{J_1}}, \dots, \theta_{X_n, \mathbf{v}_{J_n}}$.

Given a set of observations about the state s , $O(s) = \langle \mathbf{x}^{(0)}, \dots, \mathbf{x}^{(k)} \rangle$ and a new observation $\langle \mathbf{x}^{(k+1)}, s \rangle$ we have to update the values of $\theta_{\mathbf{X}, s}$ in order to maximise a combination of the current belief of the agent and the likelihood of

the entire set of observations extended with the new observation. Also in this case the agent can be more or less careful in the revision, being more or less confident in its beliefs. The update equation is therefore defined as:

$$\theta'_{\mathbf{X}, s} = \beta \cdot \theta_{\mathbf{X}, s} + (1 - \beta) \cdot \underset{\theta'}{\operatorname{argmax}} \mathcal{L}(\theta', \mathbf{x}^{(i)}, \dots, \mathbf{x}^{(k+1)}, s)$$

where the parameter $\beta \in [0, 1]$, expresses agents’ confidence in its beliefs; the higher the value of β the more careful the agent is in the revision, and

$$\mathcal{L}(\theta, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}, s) \propto \prod_{j=1}^k f(\mathbf{x}^{(j)}, s)$$

Due to the factorization of the perception function $f = \prod p_{X_i}$, we can separately update each set of parameters θ_{X_i} associated to the perception variable X_i , defining therefore

$$\theta_{X_i, s}^{k+1} = \beta \cdot \theta_{X_i, s}^k + (1 - \beta) \cdot \underset{\theta'_{X_i}}{\operatorname{argmax}} \prod_{j=1}^k p_{X_i}(x_i^{(j)} | \theta'_{X_i}) \quad (6)$$

4 Explanatory Examples

Let us now show how ALP works in Example 1. In Example 2, we first show how ALP learns new states by extending the domain of state variables and weakening constraints. In Example 3, we show how ALP can deal with highly unexpected events by adapting the planning domain.

Example 2 *Let us suppose that the robot starts with the planning domain described in Example 1.*

1. *Suppose that the robot believes to be in the state s_0 where $\text{loc}(r) = 0$, $\text{loc}(p) = 1$, and $\text{loaded} = 0$ (shortly written as $s_0 = 010$), and that the world is in the state shown in the top-left rectangle of Figure 1. Suppose that EXPLORE generates the action E (line 4) and the execution of this action moves the robot of about one unit in the east direction. The observation returned after the execution (line 5) is $\mathbf{x} = \langle x, y, t, w \rangle$ with $x \approx 1.5$, because the robot is moving east of approximately 1 unit; $y \approx 0.5$, because the robot is moving approximately horizontally; $t \approx 0$, because, differently from the model, the pack is not in that room; finally $w \approx 0$, since the robot is carrying nothing.*

2. *ALP computes the set of states $S'_0 \subseteq S$ that are above the threshold (line 6). Let us suppose that $\epsilon = 0.5$, i.e., we decide to balance our trust in the initial set of states and in the perceptions after executing actions. The robot position perception variables x and y indicate that the robot is in room 1 ($\text{loc}(r) = 1$). The sensor tag perception variable t indicates that the pack is not in the same room of the robot, i.e., $\text{loc}(p)$ is equal to 0, or 2, or 3. The weight perception variable w indicates that the robot is not loaded, i.e., $\text{loaded} = 0$. Therefore, $S'_0 = \{100, 120, 130\}$.*

3. *Since $S'_0 \neq \emptyset$, ALP computes the set S'_0 of the states in S'_0 that maximise $f(x, y, t, w, s) \cdot \text{sim}(s, 110 | \delta)$. (line 15). Notice that $\gamma(E, s_0) = \gamma(E, 010) = 110$. Notice that*

in all the states $s \in S'_0$, $f(x, y, t, w, s)$ is the same, and it approximately equal to

$$\mathcal{N}(1.5 \mid \mu = 1.5, \sigma = 1) \cdot \mathcal{N}(0.5 \mid \mu = 0.5, \sigma = 1) \cdot B(0 \mid \alpha = 1, \beta = 2) \cdot \Gamma(0 \mid k = 0, \theta = 1)$$

i.e., the robot is unloaded and in room 1, and the pack is in a different room. The values of the factor $\text{sim}(s, 110 \mid \delta)$ is also the same for all the elements of S'_0 . Therefore ALP randomly select one state of S'_0 . Suppose that ALP selects $s'_0 = 130$.

4. Then ALP jumps to line 19 and the transition $\langle 010, E, 130 \rangle$ is added to the transition log T , and the observation $\langle 130, \mathbf{x} \rangle$, with $\mathbf{x} \approx \langle 1.5, 0.5, 0, 0 \rangle$ is added to the observation log O (line 20).

5. Then ALP revises the transition function γ (line 21). According to equation (5), with $a = E$ and $s = 010$, we have:

s'	equation (5)	$\alpha = 0$	$\alpha = \frac{1}{2}$	$\alpha = 1$
130	$\alpha \cdot 0 + (1 - \alpha) \cdot 1$	1	$\frac{1}{2}$	0
110	$\alpha \cdot 1 + (1 - \alpha) \cdot 0$	0	$\frac{1}{2}$	1
others	$\alpha \cdot 0 + (1 - \alpha) \cdot 0$	0	0	0

If $\alpha > 1/2$ then γ will not be changed, otherwise $\gamma(010, E) = 130$. Let us suppose that $\alpha > 1/2$.

6. The new current state s_0 is set to 130 (line 23), and a new action is generated by EXPLORE (line 4). Let's suppose it is again E. The values returned by the perception function are $x \approx 2.5$ and $y \approx 0.5$, since the action E moves the robot east of one unit (this is possible since actually there is a room east of room 1); $t \approx 1$ (since now the pack is actually in the same room of the robot), and $w \approx 0$ (since the robot is unloaded).

7. Now there are no states in S that are above the threshold (line 6), since $p_X(2.5 \mid s)$ is very low for all the states $s \in S$. Therefore $S'_0 = \emptyset$.

8. ALP checks therefore if there are assignments to state variables that are not states in S that have the perception function above the threshold (line 8). Even in this case, for the same reason, no assignment allows for a perception function that is above the threshold. Therefore S'_0 is again empty.

9. ALP generates therefore a new state by extending the domain of state variables (line 10). EXTENDDOM starts by computing the states that maximizes the likelihood of observing $\langle x, y, t, w \rangle \approx \langle 2.5, 0.5, 1, 0 \rangle$, i.e., $s = 110$. Notice that $P_Y(\approx 0.5 \mid \text{loc}(r) = 1)$ is close to the maximum of P_Y ; similarly for $P_T(\approx 1 \mid \text{loc}(r) = \text{loc}(p))$ and $P_W(\approx 0 \mid \text{loaded} = 0)$ are also close to the maximum of P_T and P_W respectively. So if ϵ is small enough (i.e., the robot is enough "open" to the introduction of new states), X is the only variable for which $P_X(\approx 2.5 \mid \text{loc}(r) = 1)$ is below the threshold. Therefore $\mathbf{X}_{\leq \text{lik}} = \{X\}$, and $\mathbf{D}_H = \{\text{room}\}$ (line 3 of algorithm EXTENDDOM).

10. The domain room is therefore extended with a new value, obtaining $\text{room} = \{0, 1, 2, 3, 4\}$. Notice that, since room is also the domain of the variable $\text{loc}(p)$, this implies that we also extend the domain of this variable. With this extension we pass from $4 \cdot 4 \cdot 2 = 32$ possible assignments to $5 \cdot 5 \cdot 2 = 50$ possible assignments.

11. EXTENDF (line 11) extends the perception function for the new assignments: $p_X(x \mid \text{loc}(r) = 4) = \mathcal{N}(x \mid \mu_{\text{loc}(r)=4} \approx 2.5, \sigma = 1)$, and $p_Y(y \mid \text{loc}(r) = 4) = \mathcal{N}(x \mid \mu_{\text{loc}(r)=4} \approx 0.5, \sigma)$. $p_T(t \mid s)$ when s contains the new value 4 is already defined, and the perception function for W is not extended since it is not related to the state variable with domain room.

12. The s'_0 that maximises the new perception function is then 440 (lines 15 and 15). ALP therefore updates the constraints in order to include only 440 as a new state. According to Formula (4), ALP generates the following new constraints:

$$\text{loc}(r) = 4 \rightarrow \text{loc}(p) = 4 \wedge \text{loaded} = 0 \quad (7)$$

$$\text{loc}(p) = 4 \rightarrow \text{loc}(r) = 4 \wedge \text{loaded} = 0 \quad (8)$$

According to Formula (3), ALP updates the previous constraint as follows:

$$(\text{loaded} = 1 \rightarrow \text{loc}(r) = \text{loc}(p)) \vee$$

$$(\text{loc}(r) = 4 \wedge \text{loc}(p) = 4 \wedge \text{loaded} = 0)$$

which is equivalent to $\text{loaded} = 1 \rightarrow \text{loc}(r) = \text{loc}(p)$. Therefore ALP adds only the constraints (7) and (8).

13. T becomes $\langle \langle 010, E, 130 \rangle, \langle 130, E, 440 \rangle \rangle$; O becomes approximately the list of $\langle \approx \langle 0.5, 0.5, 0, 0 \rangle, 010 \rangle$, $\langle \approx \langle 1.5, 0.5, 0, 0 \rangle, 130 \rangle$, and $\langle \approx \langle 2.5, 0.5, 1, 0 \rangle, 440 \rangle$.

14. Suppose that the parameters α and β are high enough not to affect the change of γ and that they have a minimal effect on the perception function. The new state s_0 is now set to 440.

15. Suppose that EXPLORE returns action load, L. The new perceived values are approximately $\langle x, y, t, w \rangle \approx \langle 2.5, 0.5, 1, 1 \rangle$. None of the states in S is such that $p_W(w \mid s)$ is above the threshold (line 6). ALP checks therefore if there is some assignment that does not satisfy the constraints with a better likelihood (line 8). Indeed the assignment 441 is such that all the likelihoods $p_X(2.5 \mid \text{loc}(r) = 4)$, $p_Y(0.5 \mid \text{loc}(r) = 4)$, $p_T(1 \mid \text{loc}(r) = 4, \text{loc}(p) = 4)$, and $p_W(1 \mid \text{loaded} = 1)$ are above the threshold. This means that $s'_0 = 441$ is the new state, and ALP adds it to S (line 17). Adding 441 to the set of states amounts to revise the constraints following formula (3). After some simplification, ALP obtains the constraints

$$\text{loc}(r) = 4 \rightarrow \text{loc}(p) = 4 \quad (9)$$

$$\text{loc}(p) = 4 \rightarrow \text{loc}(r) = 4 \quad (10)$$

$$(\text{loaded} = 1 \rightarrow \text{loc}(r) = \text{loc}(p)) \quad (11)$$

Example 3 We continue the previous example showing how ALP adapts the planning domain to unexpected situations.

1. Suppose that EXPLORE returns the action W (go west) and that while executing this action the pack unexpectedly falls down and remains in room C, and simultaneously the cat (with a similar weight of the pack) jumps on top of the robot (!) (see bottom-right rectangle of Figure 1).

2. The sensors return $\mathbf{x} \approx \langle 1.5, 0.5, 0, 1 \rangle$. $f(\mathbf{x}, s)$ is very low (below the threshold) for all the states in S since $w \approx 1$ should mean that the pack is loaded, while $t \approx 0$ tells us that the pack is not in the same room of the robot, and the constraint (2) imposes that $\text{loc}(r) = \text{loc}(p)$.

3. ALP now checks the assignments to state variables that are not states in S (line 8). The assignment corresponding to the actual situation, i.e., the robot is loaded and in a different room from the pack, that is 141, maximises $f(\mathbf{x}, s) \cdot \text{sim}(s, \gamma(441, W) \mid \delta)$ (with $\delta < 1$). To extend S with 141 (line 17), ALP weakens the constraints according to rules (3) and (4) obtaining:

$$\begin{aligned} \text{loc}(r) = 4 &\rightarrow \text{loc}(p) = 4 \\ \text{loc}(p) = 4 &\rightarrow \text{loc}(r) = 4 \vee (\text{loc}(r) = 1 \wedge \text{loaded} = 1) \\ \text{loaded} = 1 &\rightarrow \text{loc}(r) = \text{loc}(p) \vee (\text{loc}(r) = 1 \wedge \text{loc}(p) = 4) \end{aligned}$$

4. Finally, suppose that, while the robot is carrying the pack, the cat jumps on top of the pack. The perception variable W will return a value around 2 (1 for the pack plus 1 for the cat) and $p_W(w \mid \text{loaded})$ will be below the threshold for all the states $s \in S$. Then ALP will extend the domain of the Boolean variable loaded, which becomes a three-valued variable, i.e., $\text{nr_of_carried_objects}$ is extended from $\{0, 1\}$ to $\{0, 1, 2\}$.

5 Related Work

As far as we know, the problem addressed in this paper is novel, as well as the approach and the proposed solution. Some works on domain model acquisition focus on the problem of learning action schema from collections of plans, see, e.g. [Gregory and Cresswell, 2016; McCluskey *et al.*, 2009; Cresswell *et al.*, 2013; Mourão *et al.*, 2012; Mehta *et al.*, 2011; Zhuo and Yang, 2014]. They do not consider perceptions and the set of states is given.

Works on learning and planning in POMDP (see, e.g., [Ross *et al.*, 2011; Katt *et al.*, 2017]) learn a model of the POMDP domain through interactions with the environment, with the goal to do planning, e.g., by reinforcement learning or by sampling methods. They learn the transitions, while the set of states is given as well as the mapping through observations.

Some works on POMDP Model Learning, see, e.g., [van Otterlo, 2009; Zheng *et al.*, 2018], drop the assumption that the set of states is given or the bound on the number of states is known. Two main differences with our work still exist. First, we do not learn a POMDP model, we learn a deterministic model that enables efficient planning techniques. Second, we learn the set of states represented through state variables and constraints, which is the practical way to represent a planning domain.

Our approach shares some similarities with the work on planning by reinforcement learning [Kaelbling *et al.*, 1996; Sutton and Barto, 1998; Geffner and Bonet, 2013; Yang *et al.*, 2018; Parr and Russell, 1997; Ryan, 2002; Leonetti *et al.*, 2016], since we learn by acting in the environment. However, these works focus on learning policies and assume the set of states and the correspondence between continuous data from sensors and states are fixed.

Different approaches are those followed by LatPlan and Causal InfoGAN. Causal InfoGAN [Kurutach *et al.*, 2018] learns discrete or continuous models from high dimensional sequential observations. This approach fixes a priori the size of the discrete domain model, and performs the learning off

line. Differently from our approach their goal is to generate an execution trace in the high dimensional space. LatPlan [Asai and Fukunaga, 2018] takes in input pairs of high dimensional raw data (e.g., images) corresponding to transitions. It also takes an offline approach. Our approach is online and local, we can therefore deal with a dynamic environment.

A complementary approach is pursued in works that plan and learn directly in a continuous space, see e.g., [Abbeel *et al.*, 2006; Mnih *et al.*, 2015; Co-Reyes *et al.*, 2018]. These approaches do not require a perception function, since there is no abstract discrete model of the world. Such approaches are very suited to address some tasks, e.g., moving a robot arm to a desired position or performing some manipulations. However, we believe that, in several situations, it is conceptually appropriate and practically efficient to learn an abstract discrete and deterministic model where planing is much easier and efficient to perform.

Finally, we share the idea of a planning domain at the abstract level with all the work on abstraction on MDP models, see, e.g., [Abel *et al.*, 2018]. However, our problem and approach is substantially different, since in the work on abstraction on MDP models the mapping between the original MDP and the abstract states is given, while we learn it.

6 Conclusion and Future Work

We believe this work opens a new perspective in learning planning domains and perceptions through continuous observations. The framework provides the ability to learn domains represented with state variables and constraints, which is the natural way to represent planning domains. Learning a finite deterministic planning domain represented with state variables opens up the possibility to use all the available efficient planners to reason at the abstract level. Learning the perception function takes into account the fact that, while an agent can conveniently plan at the abstract level, it perceives the world and acts through sensors and actuators that work in a continuous space. Learning perception functions allows us to learn new states that represent unexpected situations of the world. Finally, the framework allows us to learn domains incrementally, and to adapt to a changing environment.

Still a lot of work remains to do. A proof of convergence to coherent models should be provided, and the conditions of convergence should be defined. The framework should be implemented and an experimental evaluation should be performed. Additional work needs to be done to support more sophisticated action and constraint revisions on the basis of the observed transition. Finally, the ALP algorithm should be integrated with a state-of-the-art on-line planner and with efficient exploration techniques.

References

- [Abbeel *et al.*, 2006] Pieter Abbeel, Morgan Quigley, and Andrew Y. Ng. Using inaccurate models in reinforcement learning. In *ICML*, 2006.
- [Abel *et al.*, 2018] David Abel, Dilip Arumugam, Lucas Lehnert, and Michael L. Littman. State abstractions for lifelong reinforcement learning. In *ICML*, 2018.

- [Asai and Fukunaga, 2018] Masataro Asai and Alex Fukunaga. Classical planning in deep latent space: Bridging the subsymbolic-symbolic boundary. In *AAAI*, 2018.
- [Co-Reyes *et al.*, 2018] J. D. Co-Reyes, Y. Liu, A. Gupta, B. Eysenbach, P. Abbeel, and S. Levine. Self-consistent trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings. In *ICML 2018*, 2018.
- [Cresswell *et al.*, 2013] Stephen Cresswell, Thomas Leo McCluskey, and Margaret Mary West. Acquiring planning domain models using *LOCM*. *Knowledge Eng. Review*, 28(2):195–213, 2013.
- [Geffner and Bonet, 2013] H. Geffner and B. Bonet. *A Concise Introduction to Models and Methods for Automated Planning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2013.
- [Gregory and Cresswell, 2016] Peter Gregory and Stephen Cresswell. Domain model acquisition in the presence of static relations in the LOP system. In *IJCAI*, 2016.
- [Henaff *et al.*, 2017] Mikael Henaff, William F. Whitney, and Yann LeCun. Model-based planning in discrete action spaces. *CoRR*, abs/1705.07177, 2017.
- [Kaelbling *et al.*, 1996] Leslie Pack Kaelbling, Michael L. Littman, and Andrew W. Moore. Reinforcement learning: A survey. *J. Artif. Intell. Res.*, 4:237–285, 1996.
- [Katt *et al.*, 2017] Sammie Katt, Frans A. Oliehoek, and Christopher Amato. Learning in pomdps with monte carlo tree search. In *ICML*, 2017.
- [Kurutach *et al.*, 2018] Hanard Kurutach, Aviv Tamar, Ge Yang, Stuart Russell, and Pieter Abbeel. Learning plannable representations with causal infogan. In *NIPS*, 2018.
- [Leonetti *et al.*, 2016] M. Leonetti, L. Iocchi, and P. Stone. A synthesis of automated planning and reinforcement learning for efficient, robust decision-making. *Artif. Intell.*, 241:103–130, 2016.
- [McCluskey *et al.*, 2009] Thomas Leo McCluskey, Stephen Cresswell, N. Elisabeth Richardson, and Margaret Mary West. Automated acquisition of action knowledge. In *ICAART*, 2009.
- [Mehta *et al.*, 2011] Neville Mehta, Prasad Tadepalli, and Alan Fern. Autonomous learning of action models for planning. In *NIPS*, 2011.
- [Mnih *et al.*, 2015] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [Mourão *et al.*, 2012] Kira Mourão, Luke S. Zettlemoyer, Ronald P. A. Petrick, and Mark Steedman. Learning STRIPS operators from noisy and incomplete observations. In *UAI*, 2012.
- [Parr and Russell, 1997] R. Parr and S. J. Russell. Reinforcement learning with hierarchies of machines. In *NIPS*, 1997.
- [Pasula *et al.*, 2004] Hanna Pasula, Luke S. Zettlemoyer, and Leslie Pack Kaelbling. Learning probabilistic relational planning rules. In *ICAPS*, 2004.
- [Pasula *et al.*, 2007] Hanna M. Pasula, Luke S. Zettlemoyer, and Leslie Pack Kaelbling. Learning symbolic models of stochastic domains. *J. Artif. Intell. Res.*, 29:309–352, 2007.
- [Ross *et al.*, 2011] Stéphane Ross, Joelle Pineau, Brahim Chaib-draa, and Pierre Kreitmann. A bayesian approach for learning and planning in partially observable markov decision processes. *Journal of Machine Learning Research*, 12:1729–1770, 2011.
- [Ryan, 2002] M. R. K. Ryan. Using abstract models of behaviours to automatically generate reinforcement learning hierarchies. In *ICML*, 2002.
- [Sutton and Barto, 1998] R. S. Sutton and A. G. Barto. *Reinforcement learning - an introduction*. Adaptive computation and machine learning. MIT Press, 1998.
- [van Otterlo, 2009] Martijn van Otterlo. *The Logic of Adaptive Behavior - Knowledge Representation and Algorithms for Adaptive Sequential Decision Making under Uncertainty in First-Order and Relational Domains*, volume 192 of *Frontiers in Artificial Intelligence and Applications*. IOS Press, 2009.
- [Yang *et al.*, 2007] Q Yang, K Wu, and Y Jang. Learning action models from plan examples using weighted maxsat. *Artif. Intell.*, 171:107–143, 2007.
- [Yang *et al.*, 2018] F. Yang, D. Lyu, B. Liu, and S. Gustafson. PEORL: integrating symbolic planning and hierarchical reinforcement learning for robust decision-making. In *IJCAI*, 2018.
- [Zettlemoyer *et al.*, 2005] Luke S. Zettlemoyer, Hanna Pasula, and Leslie Pack Kaelbling. Learning planning rules in noisy stochastic worlds. In *AAAI*, 2005.
- [Zheng *et al.*, 2018] Wei Zheng, Bo Wu, and Hai Lin. POMDP model learning for human robot collaboration. In *57th IEEE Conference on Decision and Control, CDC 2018, Miami, FL, USA, December 17-19, 2018*, pages 1156–1161, 2018.
- [Zhuo and Kambhampati, 2013] Hankz Hankui Zhuo and Subbarao Kambhampati. Action-model acquisition from noisy plan traces. In *IJCAI*, 2013.
- [Zhuo and Yang, 2014] Hankz Hankui Zhuo and Qiang Yang. Action-model acquisition for planning via transfer learning. *Artif. Intell.*, 212:80–103, 2014.
- [Zhuo *et al.*, 2010] Hankz Hankui Zhuo, Qiang Yang, Derek Hao Hu, and Lei Li. Learning complex action models with quantifiers and logical implications. *Artif. Intell.*, 174(18):1540–1569, 2010.