

LUND UNIVERSITY

Learning and Planning of Situated Resource Bounded Agents

Nowaczyk, Slawomir

Published in:

Proceedings of the 23rd Annual Workshop of the Swedish Artificial Intelligence Society

2006

Link to publication

Citation for published version (APA): Nowaczyk, S. (2006). Learning and Planning of Situated Resource Bounded Agents. In Proceedings of the 23rd Annual Workshop of the Swedish Artificial Intelligence Society

Total number of authors: 1

General rights

Unless other specific re-use rights are stated the following general rights apply: Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

• Users may download and print one copy of any publication from the public portal for the purpose of private study

or research.
You may not further distribute the material or use it for any profit-making activity or commercial gain

· You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: https://creativecommons.org/licenses/

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

LUND UNIVERSITY

PO Box 117 221 00 Lund +46 46-222 00 00

Learning and Planning of Situated Resource Bounded Agents*

Sławomir Nowaczyk Slawomir.Nowaczyk@cs.lth.se Department of Computer Science Lund University, Sweden

Abstract

This paper presents an investigation of rational agents that have limited computational resources and that can interact with their environments. We analyse how such agents can combine deductive reasoning using domain knowledge and inductive learning from past experiences, while remaining time-aware in a manner appropriate for beings situated in a dynamic universe. In particular, we consider how they can create and reason about partial plans, choose and execute the best ones of them — in such way as to acquire the most knowledge. We also discuss what are the different types of interactions with the world and how they can influence agent's ability to consciously direct its own learning process.

1 Introduction

In our research we are interested in building rational agents that can interact with their environment. In order to be practically useful, such agents should be modelled as having bounded computational resources. Moreover, since they are situated in a dynamic world, they need to be aware of the notion of time — in particular, that their reasoning process is not instantaneous. On the other hand, such agents have the possibility to acquire important knowledge by observing the environment surrounding them and by analysing their past interactions with it.

This paper focuses on discussion how reasoning machinery of such agents can adapt to various models of interaction with the world. We also present how

^{*}This is an extended version of paper submitted to Student session on ESSLLI 2006, combined with some ideas presented at IJCAI 2005 Workshop on Planning and Learning in A Priori Unknown or Dynamic Domains

such rational agents can deal with planning in domains where complexity makes finding complete solutions intractable. Clearly, it is often not realistic to expect an agent to be able to find a total plan which solves a problem at hand. Therefore, we investigate how an agent can create and reason about *partial plans*. By that we mean plans which bring it somewhat closer to achieving the goal, while still being simple and short enough to be computable in reasonable time. Currently we mainly focus on plans which allow an agent to acquire additional knowledge about the world.

By executing such "information-providing" partial plans, an agent can greatly simplify subsequent planning process — it no longer needs to take into account the vast number of possible situations which will be inconsistent with newly observed state of the world. Thus, it can proceed further in a more effective way, by devoting its computational resources to more relevant issues.

If the environment is modelled sufficiently well (for example, if a simulator exists), the agent may have a high degree of freedom in exploring it and in deciding how to interact with it. It may be possible to gain information that the agent would not be able to, by itself, observe directly. In many domains it is significantly easier to build and employ a simulator than to analytically predict results of complex interactions. In other cases, for example when the agent is a robot situated in an unknown environment, it must learn "in the wild" and be aware that the actions it executes are final: they do happen and there is no way of undoing them, other than performing, if possible, a reverse action.

In order to accommodate all of the above we use a variant of Active Logic (Elgot-Drapkin *et al.* 1999) as agent's reasoning formalism. It was designed for non-omniscient agents and has mechanisms for dealing with uncertain and contradictory knowledge. We believe that Active Logic is a good reasoning technique for versatile agents, as it has been successfully applied to several different problems, including some in which planning plays a very prominent role (Purang *et al.* 1999). Moreover, in order to be able to intentionally direct its own learning process, the agent needs to reason about its own knowledge and lack of it — thus, its logic needs to be augmented with epistemic concepts (Fagin *et al.* 1995).

In other words, our agents are supposed to combine deductive and inductive reasoning with time-awareness. We believe that the interactions among those three aspects are crucial for developing truly intelligent systems. It is not our goal to analyse strict deadlines or precise time measurements (although we do not exclude a possibility of doing that), but rather to express that a rational agent needs the ability to reason about committing its resources to various tasks (Chong *et al.* 2002).

2 Wumpus Game

The example problem we will be using through this paper is a well-known game of Wumpus, a common testbed for intelligent agents. In its basic form, the game takes place on a board through which an player is allowed to move freely. A beast called Wumpus occupies one, initially unknown, square. Agent's goal is to kill the creature, a task that can be achieved by shooting an arrow on that square. Luckily, Wumpus is a smelly creature, so the player always knows if the monster is nearby. Unfortunately, not in which direction, exactly. At the same time, when walking around, the player can get eaten by the monster if he stumbles across it.

This game is concise enough to be explained easily, but finding a solution is sufficiently complex to illustrate the issues we want to emphasise. We look at it as one instance of a significantly broader class of problems, along the lines of *General Game Playing*, where an agent accepts a formal description of an arbitrary game and, without further human interaction, can play it effectively.

3 Agent Architecture

We use a simple architecture for our agent, as presented in Fig. 1. It consists of three main elements, corresponding to the three main tasks of the agent.

The Deductor reasons about world, possible actions and what could be their consequences. Its main aim is to generate plans applicable in current situation and predict — at least as far as past experience, imperfect domain knowledge and limited computational resources allow — effects each of them will have, in particular what new knowledge can be acquired.

The Actor is responsible for overseeing the reasoning process, mainly for introducing new observations into the knowledge base and for choosing plans for execution. Basically, it decides *when* to switch from deliberation to acting, and which of the plans under consideration to execute.

These two modules form the core of the agent. By creating and executing a sequence of partial plans our agent moves progressively closer and closer to its



Figure 1: The architecture of the system.

goal, until it reaches a point where a solution can be directly created by Deductor.

The learning module is necessary in order to ensure that the plans agent chooses for execution are indeed "good" ones. After the game is over, regardless of whether the agent has won or lost, learning system inductively generalises experience it has gathered — attempting to improve Deductor's and Actor's performance. Our goal is to use the learned information to fill gaps in the domain knowledge, to figure out generally interesting reasoning directions, to discover relevant subgoals and, finally, to more efficiently select the best partial plan.

4 Logical Reasoning

The language used by Deductor is the First Order Logic augmented with Situation Calculus mechanisms. Within a given situation, knowledge is expressed using standard FOL. In particular, we do not put any limitations on the expressiveness of the language. Predicate Knows describes knowledge of the agent, e.g.,

$$Knows[smell(a) \leftrightarrow \exists_x(Wumpus(x) \land Neighbour(a, x))]$$

means: agent knows that it smells on exactly those squares which neighbour Wumpus' position. The predicate Knows may be nested, although it is useful only in a couple of specialised contexts. We use standard reification mechanism for putting formulae as parameters of a predicate (Reiter 2001).

In order to describe actions and changes, we employ a well-known Situation Calculus approach, using a predicate *Holds* (*situation, formula*) to denote that the *formula* holds in *situation* and a predicate *Informs* (*action, groundedwff*) to denote that *action* provides information whether *groundedwff* holds. We also introduce function *Result* (*situation, action*), which returns the set of situations resulting from applying *action* in *situation*.

In order to facilitate agent's reasoning about changing world, we treat predicate Knows in a similar way as Holds, i.e. we introduce an additional parameter denoting the current situation. Moreover, since the agent needs to reason about knowledge-producing actions, we add yet another parameter, namely the plan agent is going to execute. Therefore, the previous formula should actually be written as:

$$Knows[s, p, smell(a) \leftrightarrow \exists_x(Wumpus(x) \land Neighbour(a, x))]$$

and mean: agent knows that executing plan p in situation s leads to a new situation, such that it smells on exactly those squares which neighbour Wumpus' position. This particular formula is an universal law of the world, valid regardless of the chosen s and p, but many interesting ones — e.g. "Wumpus(a)" or "Knows[smell(b)]" — are true only for specific s and p. As we mentioned earlier, the agent employs Active Logic — a formalism intended to describe the deduction as an ongoing process, as opposed to characterising just some static, fixed-point consequence relation. To this end, it annotates every formula with a time-stamp (usually an integer) of when it was first derived, incrementing the label with every application of an inference rule:

$$\frac{i: a, a \to b}{i+1: b}$$

Additional features which are available in AL and important for this work include the Now predicate, true only during current time point (i.e., "i : Now(j)" is true for all i = j, but false for all $i \neq j$) and the observation function, delivering axioms that are valid since a specific point in time. They can be used to naturally model agent acquiring new knowledge from the environment, including changes which are external to the agent. This way AL lifts two important limitations of the classical Situation Calculus. Finally, the logic has provisions for dealing with contradictory knowledge, useful when agent learns something which is not completely correct and that later conflicts with incoming observation.

The predicate *Now* makes it possible to reason about passing time and, combined with observation function delivering knowledge about external events, allows the agent to remain responsive during its deliberations. This way both Deductor and Actor can keep track of how the reasoning is progressing and make informed decisions about balancing thinking and acting.

The plans agent reasons about consist of a concatenation of classical and conditional actions, the latter of the form (*predicate* ? $action_1 : action_2$), with the usual meaning that $action_1$ will be executed if *predicate* holds, and $action_2$ will be executed otherwise. For a well-developed discussion of other possible ways of representing conditional partial plans and of interleaving planning and execution see, for example, Bertoli, Cimatti, & Traverso (2004)

One of the reasons we have chosen symbolic representation of plans, as opposed to a policy (an assignment of value to each state–action pair) is that we intend to deal with other types of goals than just reachability ones. For a discussion of possibilities and rationalisation of why such goals are interesting, see for example Bertoli *et al.* (2003), where authors present a solution for planning with goals described in Computational Tree Logic. This formalism allows to express goals of the kind "value of *a* will never be changed", "*a* will be eventually restored to its original value" or "value of *a*, after time *t*, will always be *b*" etc.

To summarise, our agent uses AL to reason about its own knowledge, which is very important in the Wumpus domain. Here, the main goal can be reduced to "learn the position of Wumpus", so active planning for knowledge acquisition is crucial. Agent also requires an ability to compare what kind of information will execution of each plan provide, in order to be able to choose the best one of them.

5 Actor

The Actor module supervises the deduction process and breaks it at selected moments, e.g., when it notices a particularly interesting plan or when it decides that sufficiently long time has been spent on planning. It then *evaluates* existing partial plans and executes the best one of them. The evaluation process is crucial here, and we expect the subsequent learning process to greatly contribute to its improvement. In the beginning, the choice may be done at random, or some simple heuristic may be used. After execution of partial plan, a new situation is reached and the Actor lets the Deductor create another set of possible plans.

This is repeated as many times as needed, until the game episode is either won or lost. Losing the game clearly identifies bad choices on the part of the Actor and leads to an update of the evaluation function.

Winning the game also yields feedback that may be used for improving this function, but it also provides a possibility to (re)construct a complete plan, i.e. one which starts in the initial situation and ends in a winning state. If such a plan can be found, it may be subsequently used to quickly solve any problem instance for which it is applicable. Moreover, even if such plan is not directly applicable, an Actor can use it when evaluating other plans found by the Deductor. Those with structure similar to the successful one are more likely to be worthwhile.

6 Learning

When analysing learning module, it is important to keep in mind that our agent has a dual aim, akin to the exploration and exploitation dilemma in reinforcement learning. On one hand, it wants to win the current game episode, but at the same time it needs to learn as much general knowledge as possible, in order to improve its future performance.

Currently we are mainly investigating the learning module from Actor's perspective — using ILP to evaluate quality of partial plans is, to the best of our knowledge, a novel idea. One issue is that work on ILP has been dealing almost exclusively with the problem of *classification*, while our situation requires *evaluation*. There is no predefined set of classes into which plans should be assigned. What our agent needs is a way to choose the *best* one of them.

For now, however, we focus on distinguishing a special class of "bad" plans, namely ones that lead to losing the game. Clearly some plans — those that in agent's experience *did* so — are bad ones. But not every plan which does not cause the agent to lose is a *good* plan. Further, not every plan that leads to *winning* a game is a good one. An agent might have executed a dangerous plan and win only because it has been lucky.

Therefore, we define as positive examples those plans which lead, or can be proven to *possibly* lead, to the defeat. On the other hand, those plans which can

be proven to *never* cause defeat are negative examples. There is a third class of plans, when neither of the above assertions can be proven. We are working on how to use such examples in learning most effectively.

Nevertheless, this is only the beginning. After all, in many situations a more "proactive" approach than simple *not-losing* is required. One promising idea is to explore the epistemic quality of plans: an agent should pursue those which provide the most important knowledge. Another way of expressing distinction between good and bad partial plans, one we feel can give very good results, is discovering relevant subgoals and landmarks, as in Hoffmann, Porteous, & Sebastia (2004).

7 Environment Interaction

One of the main contributions of our research lies in the "consciousness" of interactions between an agent and its environment, conducted in such a way as to maximise the knowledge that can be obtained. In particular, an agent is facing, at all times, the exploration versus exploitation dilemma, i.e., it both needs to gather new knowledge *and* to win the current game episode.

In order to facilitate such reasoning, our agent requires an ability to both act in the world and to observe it. Finally, it needs to consider its own knowledge and how it will (or *can*) change in response to various events taking place in the environment. In different domains and applications different models of interactions with the world are possible.

The most unrestrictive case is a simulator, where an agent has complete control over the (training) environment. It can setup an arbitrary situation, execute some actions and observe the results. Such a scenario is common in, for example, physical modelling, where it is often much easier to simulate things than to predict their behaviour and interactions. In a similar spirit, it may be easier for our agent to "ask the environment" about validity of some formula than to prove it.

If agent's freedom is slightly more restricted, it is possible that it is not allowed to freely change the environment, but can "try out" several plans in a given situation. For example, the agent may provide a set of plans and receive an outcome for each of them. Alternatively, it may store some opaque *situation identifier* so that it can revisit the same situation at later time. This model is also suitable for agents that do not have perfect knowledge of the world, as the "replay" capability does not assume *the agent* is able to fully reconstruct the situation or knows the state of the world completely.

In our opinion, this is the most interesting setting: it gives the agent sufficient freedom to allow it to achieve interesting results and at the same time is not, in many domains, overly infeasible. On the other hand, we are working on ways in which this setting could be made even more practical — one idea is having an agent accept the fact that in several replays "the same" situation could vary slightly. For example, physical agent might request an operator to restore previous

state of the world: it would not really be identical, but it may be sufficiently close. Alternatively, in some applications, only a subset of situations may be "replayable" — only those, for example, that an agent can restore, with required tolerance, all by itself.

In most applications, however, the agent is only able to influence its own actions and have no control whatsoever over the rest of the world. This is also the most suitable model for an *autonomous* physical agent. In such case, the environment will irreversibly move into the subsequent state upon each agent's action (or any other event), leaving it no option but to adapt. It may still be interesting, in some situations, to substitute acting for reasoning, but the agent needs to be aware that once acted upon, the current situation will be gone, possibly forever. It thus needs to consider if saving some deduction effort is indeed the best possible course of action, or if doing something else instead would be more advantageous.

Finally, we can imagine a physical agent situated in a *dangerous* environment, where it is not even plausible for it to freely choose its actions — it needs to, first, assert that an action is reasonably safe. In this case, unlike the previous one, a significant amount of reasoning *needs* to be performed before every experiment.

As an orthogonal issue, sometimes it is feasible for an agent to execute an action, observe the results, reason about them and figure out the next action to perform. But in many applications the "value" of time varies significantly. There are situations where an agent may freely spend its time meditating, and there are situations where decisions must be made quickly. For example, in RoboCup robotic soccer domain, when the ball is in possession of a friendly player, the agent just needs to position itself in a good way for a possible pass — a task which is not too demanding and leaves agent free to ponder more "philosophical" issues. On the other hand, when the ball is rolling in agent's direction, time is of essence and an agent better had plans ready for several most plausible action outcomes.

8 Related Work

Combination of planning and learning is an area of active research, in addition to the extensive amount of work being done separately in those respective fields.

There has been significant amount of work done in learning about what actions to take in a particular situation. One notable example is Khardon (1999), where author showed important theoretical results about PAC-learnability of action strategies in various models. In Moyle (2002) author discussed a more practical approach to learning Event Calculus programs using Theory Completion. He used extraction-case abduction and the ALECTO system in order to simultaneously learn two mutually related predicates (*Initiates* and *Terminates*) from positive-only observations. Recently, Könik & Laird (2004) developed a system which is able to learn low-level actions and plans from goal hierarchies and action examples provided by experts, within the SOAR architecture. The work mentioned above focuses primarily on learning how to act, without focusing on reaching conclusions in a deductive way. In a sense, the results are somewhat more similar to the reactive-like behaviour than to classical planning system, with important similarities to the reinforcement learning and related techniques.

One attempt to escape the trap of large search space has been presented in Džeroski, Raedt, & Driessens (2001), where relational abstractions are used to substantially reduce cardinality of search space. Still, this new space is subjected to reinforcement learning, not to a symbolic planning system. A conceptually similar idea, but where relational representation is actually being learned via behaviour cloning techniques, is presented in Morales (2004).

Recently, Colton & Muggleton (2003) showed several ideas about how to learn interesting facts about the world, as opposed to learning a description of a predefined concept. A somewhat similar result, more specifically related to planning, has been presented in (Fern, Yoon, & Givan 2004), where the system learns domain-dependent control knowledge beneficial in planning tasks.

Yet another track of research focuses on (deductive) planning, taking into account incompleteness of agent's knowledge and uncertainty about the world. Conditional plans, generalised policies, conformant plans and universal plans are the terms used (Bertoli, Cimatti, & Traverso 2004).

9 Conclusions

The work presented here is a discussion of an interesting track of research, rather than a report on some concrete results. We have introduced an agent architecture facilitating resource-aware deductive planning interwoven with plan execution and supported by inductive, life-long learning. The particular deduction mechanism used is based on Active Logic, in order to incorporate time-awareness into the reasoning itself. The plans created in deductive way are conditional, accounting for possible results of future actions, in particular information-gathering ones.

We intend to continue this work in several directions. Discovering subgoals and subplans seems to be one of the most useful capabilities of human problem solving and we would like our agent to invent and use such concept. In our example domain a useful subgoal could be "First, find a place where it smells." In addition, Deductor should be able to conceive general rules of rational behaviour, such as "Don't shoot if you don't know Wumpus' position". Yet another clear advantage would be the ability to reuse a previously successful plan in a different situation. Finally, domain experts often are an invaluable source of knowledge that the agent should be able to exploit, if possible.

The ideas above do not cover all the possible further investigations and extensions of the proposed system; it is just a biased presentation of the authors' own interests and judgements.

References

- Bertoli, P.; Cimatti, A.; Pistore, M.; and Traverso, P. 2003. A framework for planning with extended goals under partial observability. In *International Conference on Automated Planning and Scheduling*, 215–225.
- Bertoli, P.; Cimatti, A.; and Traverso, P. 2004. Interleaving execution and planning for nondeterministic, partially observable domains. In *European Conference on Artificial Intelligence*, 657–661.
- Chong, W.; O'Donovan-Anderson, M.; Okamoto, Y.; and Perlis, D. 2002. Seven days in the life of a robotic agent. In *GSFC/JPL Workshop on Radical Agent Concepts*.
- Colton, S., and Muggleton, S. 2003. ILP for mathematical discovery. In *13th International Conference on Inductive Logic Programming*.
- Džeroski, S.; Raedt, L. D.; and Driessens, K. 2001. Relational reinforcement learning. *Machine Learning* 43(1/2):7–52.
- Elgot-Drapkin, J.; Kraus, S.; Miller, M.; Nirkhe, M.; and Perlis, D. 1999. Active logics: A unified formal approach to episodic reasoning. Technical Report CS-TR-4072, University of Maryland.
- Fagin, R.; Halpern, J. Y.; Vardi, M. Y.; and Moses, Y. 1995. *Reasoning about knowledge*. MIT Press.
- Fern, A.; Yoon, S.; and Givan, R. 2004. Learning domain-specific control knowledge from random walks. In *International Conference on Automated Planning and Scheduling*.
- Hoffmann, J.; Porteous, J.; and Sebastia, L. 2004. Ordered landmarks in planning. *Journal* of Artificial Intelligence Research 22:215–278.
- Khardon, R. 1999. Learning to take actions. Machine Learning 35:57-90.
- Könik, T., and Laird, J. 2004. Learning goal hierarchies from structured observations and expert annotations. In *14th International Conference on Inductive Logic Programming*.
- Morales, E. P. 2004. Relational state abstraction for reinforcement learning. In *Proceedings of the ICML'04 Workshop on Relational Reinforcement Learning*.
- Moyle, S. 2002. Using theory completion to learn a robot navigation control program. In *12th International Conference on Inductive Logic Programming*.
- Purang, K.; Purushothaman, D.; Traum, D.; Andersen, C.; and Perlis, D. 1999. Practical reasoning and plan execution with active logic. In *Proceedings of the IJCAI-99 Workshop on Practical Reasoning and Rationality*, 30–38.
- Reiter, R. 2001. *Knowledge in Action: Logical Foundations for Specifying and Implementing Dynamical Systems.* The MIT Press.