

Time Complexity of the Incremental-Pruning Algorithm for Markov Decision Networks

H.H.L.M. Donkers J.W.H.M. Uiterwijk H.J. van den Herik
{donkers,uiterwijk,herik}@cs.unimaas.nl
Department of Computer Science, *matriks-cs*
Universiteit Maastricht
P.O. Box 616, 6200 MD Maastricht, The Netherlands

Abstract

A new representation technique for the planning of decisions under uncertainty is presented. The technique, called Markov decision networks (mdns), solves planning tasks that cannot or can hardly be solved by any other known representation technique, such as decision trees, influence diagrams, and pomdps.

We show that linear mdns are pomdps. Hence, the incremental-pruning algorithm for pomdps is applicable to this subclass of mdns. It turns out that the time complexity of the algorithm is lower for mdns than for general pomdps, due to the division of the state space into state groups. An example shows the practical benefit of our findings: an optimal strategy can be computed in a few seconds.

Contents

1	Introduction	2
2	Markov decision networks	3
2.1	A framework for mdns	3
2.2	A graphical representation	4
2.3	Semantics	4
3	Solving MDNs for restricted strategy sets	5
3.1	Decision-tree approach	5
3.2	History rules	5
4	Solving linear MDNs	6
4.1	Linear mdns as pomdps	6
4.2	Incremental pruning	7
4.3	Specializing incremental pruning for mdns	7
4.4	Linking IP-scheme results	8
4.5	The incremental-pruning algorithm for MDNs	9
5	Time complexity of the incremental-pruning algorithm for MDNs	10
6	Example	12
7	Conclusion	13
8	References	13

1. Introduction

In Artificial Intelligence, the planning of decisions for agents in uncertain environments is a main topic. The treatment of the subject is diverse. Depending on the kinds of decisions, agents, and uncertainty, a whole range of representation techniques is available to model and solve the planning tasks [8]. In this paper we concentrate on those planning problems for which (1) the uncertainty is modelled by Bayesian probabilistics, (2) there is one single agent that must build an optimal decision plan, (3) all decision options are completely known to the agent, and (4) the quality of a plan is measured only by its expected pay-off.

At present, there are several techniques available to deal with the above type of decision problem, but none of them is capable of solving all decision problems so typed satisfactorily. The oldest technique is the classical decision tree. It can only be used in small instances for which an exhaustive search through all possible plans is feasible. The influence diagram, in its original form, has as drawback that the ordering of the decision options is fixed by the structure of the diagram, and that repetition of decisions is inhibited. There are some alterations of the influence diagrams that try to alleviate these shortcomings, but full flexibility is still not reached. The technique of partially observed Markov decision processes (pomdps) offers the highest possible flexibility of combining decision options in any order. It continues the probabilities involved to behave Markovian, but it lacks the clear semantics of the influence diagrams. Moreover, even the most efficient exact solution method for pomdps (i.e., incremental pruning) is only feasible for small instances [10].

In this paper we introduce a new technique for solving a special type of decision problem: the Markov decision network (mdn)¹. As the name indicates, we use a Markovian probability model, which gives the technique the same kind of planning flexibility as the pomdp technique does. However, we combine this flexibility with a graphical representation and under-

standable semantics, both improving the capability to model decision problems. The demands imposed on the decision problems to be modelled by an mdn are severe. But once these demands have been complied with, mdns offer a clear representation and a set of efficient solution methods.

After introducing the framework for mdns and two general solution methods, we show that linear mdns, a subclass of mdns, are pomdps. This means that for these mdns, solution methods for pomdps are applicable. Especially the incremental-pruning algorithm seems to be appropriate. The particular structure of mdns enables us to improve this algorithm for mdns.

The time complexity of this improved algorithm is investigated and we show it to be lower than that of the general algorithm for pomdps. To conclude this paper, we give a small example of an mdn and its solution obtained by incremental pruning.

¹Part of this paper has been submitted for publication elsewhere ([5]).

2. Markov decision networks

The concept of mdns originates from the same source as the influence diagram, namely from the decision tree. The observation that different end nodes in a decision tree often represent the same final states of the world and that sometimes even whole subtrees are repeated in a decision tree, led Scott Olmsted in 1983 to the introduction of the concept of coalescence in decision trees [7]. An analogous observation in decision trees for medical use resulted in the introduction of the concept of mdns [2]. The idea of an mdn representation also stems from the first author's observation at the Department of Medical Informatics that medical students often have difficulties to understand and set up decision trees. Both observations motivated our group to simplify the concept of decision trees in the direction of compressing the tree into a network [2]. The idea was to provide the students with a straightforward technique to model and solve (medical) decision tasks. Of course, we were aware that such a simplification would reduce the class of decision tasks that could be handled. Later it transpired that the new concept was pivotal for solving planning tasks automatically, which is our present focus of research. Below we present a concise formal framework for mdns. The complete framework, along with the necessary proofs, is given in [3].

2.1. A framework for MDNs

An mdn is a quintuple of sets $\langle A; E; Z; W; G_i \rangle$, where A is the set of available actions $a_1; a_2; \dots$, and E the set of available experiments $e_1; e_2; \dots$. Actions change the state of the world, experiments give uncertain information on the actual state of the world by producing some outcome. For every experiment e_i , a non-empty set Z_{e_i} of outcomes (or results) $z_{e_i,1}; z_{e_i,2}; \dots$ exists. The set Z is the union of these outcome sets Z_{e_i} . The results of an action or an experiment only depend on the actual state of the world. This property, the Markov property, gave its name to the Markov decision network.

The states of the world are represented by the finite

set $W = \{w_1; w_2; \dots; w_g\}$. The set $G = \{G_1; G_2; \dots; G_g\}$ is the set of state groups, where $G_i \subseteq W$. These state groups form a partition of W . Two special state groups are the start group and the end group, represented by G_0 and G_u ¹, respectively. The start group G_0 is the group in which the world is at the start of a plan, the end group G_u is the group that must always be reached when the plan ends and in which the final rewards (utilities) are received. It is possible to have a state group that is both start group and end group.

Three functions G^{in} , G^{out} , and G^{test} determine the structure of an mdn: function $G^{in}(a)$ gives the state group for which action a is meaningful, function $G^{out}(a)$ gives the state group that has been reached after action a has been performed, and function $G^{test}(e)$ gives the state group for which experiment e is meaningful. The inverse functions, $A(G)$ and $E(G)$, give the actions and experiments that are available for group G .

The probabilities in an mdn are defined by: (1) for each action a in A : a $|G^{out}(a)| \times |G^{in}(a)|$ transition probability matrix A_a , where $A_a(i; j)$ is the probability that state i in group $G^{out}(a)$ will be reached if the world was in state j of group $G^{in}(a)$ just before the action was performed; (2) for each experiment e in E : a $|Z_e| \times |G^{test}(e)|$ sensitivity matrix E_e , where $E_e(i; j)$ is the probability that the i -th outcome of Z_e is observed by the agent when the world is in state j of group $G^{test}(e)$; and (3) for the prior probabilities: a probability vector x_0 , a $|G_0|$ -sized vector in which $x_0(i)$ denotes the prior probability of the world being in state i of group G_0 :

As a result of these probabilities, the state groups are mutually exclusive with respect to the probability distribution over the states. If there is a positive probability of the world being in one state of a group, then there is zero probability of the world being in a state of another group.

The costs in an mdn are given by the two functions $K(a)$, the cost of doing action a , and $K(e)$, the cost of

¹The subscript u in G_u refers to the fact that the states of the end group are connected to utilities.

performing experiment e . Utilities are expressed in a $jG_{u,j}$ -sized utility vector u in which $u(i) \in \mathbb{R}$ denotes the utility of the world reaching final state i of end group G_u .

Decision plans are called strategies in mdns. Strategies can be represented by strings of action symbols a_i , experiment symbols e_i , and a special ending symbol $\$$. For instance, the strategy $s = \langle e_1; a_1; \$; a_2; \$ \rangle$ denotes the decision plan to start with experiment e_1 and to perform action a_1 if the first outcome of e_1 is observed, and a_2 if the second outcome of e_1 is observed. Strategies that can actually be realized in an mdn are called legal.

An episode is a sequence of events (actions and experiments) that could actually happen if the agent follows a legal strategy. Episodes can be represented by simple strings of symbols for actions and outcomes of experiments. (The string $h = \langle e_1; 1; a_1 \rangle$ denotes the episode of experiment e_1 resulting in the first outcome, whereafter action a_1 is performed).

The (pay-off) value $V(h)$ of an episode h is some function of the collected costs $K(h)$ and the expected utility $U(h)$ of the episode. The form of this value function $V(h)$ can be chosen at will. However, any choice does influence (or even determine) the solution method to be used for solving the mdn. The expected value $V(s)$ of a strategy s is the expected (pay-off) value of the episodes that can happen if that strategy is followed:

$$V(s) = \sum_{h \in s} P(h)V(h) \quad (2.1)$$

In this formula, $P(h)$ is the probability that episode h occurs if strategy s is followed. Finally, the general task of solving an mdn is to find a legal strategy s^* which has the highest expected value $V(s^*)$ among all legal strategies.

2.2. A graphical representation

The graphical representation of an mdn is an augmented directed graph. The graph has a node for every state group, represented by an oval. The names of the states in the group can be written inside the oval. If needed, the groups themselves can also be named. This group name is placed outside the oval. The start group and end group can be marked specifically. Actions are represented by directed arcs between the nodes. The names of the actions are writ-

ten alongside the arcs. Finally, experiments are represented as triangles that are linked to the oval of their test group. The outcomes of an experiment are not represented in the graph.

Figure 2.1 gives an example of the graphical representation for an mdn that models a simple diagnostic planning task for repairing parts. In this case there are two state groups, each containing two states. There are also two actions ("Repair" and "Don't Repair") and two experiments ("test 1" and "test 2"). The planning task is to decide how to combine both tests in such a manner that a maximum number of parts is fixed while keeping the total costs as low as possible. On the base of this example we briefly discuss some extra semantics that mdns give rise to, in comparison with influence diagrams and pomdps.

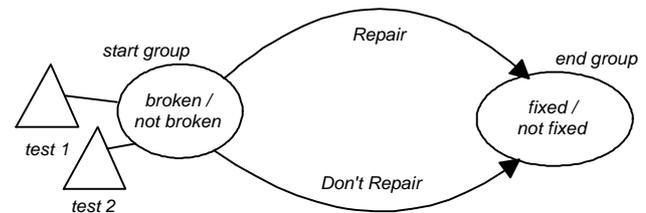


Figure 2.1: The graphical representation of an mdn.

2.3. Semantics

The meaning of a state group in mdns is broader than just an arbitrary set of states. State groups often denote specific views on the world that the agent uses at different moments in time. In the example of Figure 2.1, the first state group represents the view in which parts are either broken or not broken. The second group represents the view in which the agent looks upon parts as being fixed or not fixed. According to the graph, performing diagnostic tests is not meaningful in this stage and neither is reparation. Moving from one state group to another means not only changing the world, but also changing the perspective of the agent. Sometimes, the change of perspective is even the only alteration that is caused by an action. Hence, an action in an mdn means more than just having an effect on the environment: it also changes, and sometimes exclusively, the inner state of the agent.

3. Solving MDNs for restricted strategy sets

MdNs were developed not only to model planning tasks, but also to provide a way of solving these planning tasks automatically. In this section, two solution methods for mdNs are discussed that restrict the set of investigated strategies, viz. the standard decision-tree approach, and the decision-tree approach enriched with history rules.

3.1. Decision-tree approach

The decision-tree solution method is straightforward but restricted. It can be useful in many practical cases. In this approach, the mdN is used only as some kind of language through which a decision tree is constructed by the agent. The decision tree contains all plans that the agent judges to be relevant to investigate. The available decision options for the construction of the decision tree are derived from the arcs in the network. The branches of the chance nodes are constructed in accordance with the outcomes of the experiments.

After the construction of the tree, the quantitative information of the mdN is used to fill in the probabilities and utilities of the decision tree automatically, and the classical folding-back procedure is used to solve the decision tree. In [2] a computer program is described that offers this solution method in a user-friendly way.

3.2. History rules

The second solution method is an extension of the decision-tree approach. Instead of specifying a complete decision tree, the agent just states a set of rules that the generated plans should comply with. These rules are called history rules, because they should be able to decide whether a decision option is feasible given a history of actions and outcomes of experiments. Examples of history rules are: "never repeat an action" and "allow this action only after a positive test outcome." In [3, 4] an algorithm is presented that solves an mdN, given a set of history rules. It is imposed that these rules may not allow infinite episodes, since then the algorithm will not end.

The history-rule algorithm handles four tasks in one run: it collects costs, computes probabilities, calculates expected values, and constructs an optimal strategy. In essence, it is the same as performing the folding-back procedure. The main difference is that the decision tree itself is never constructed explicitly.

The complexity of this algorithm depends primarily on the history rules given. For instance, if one history rule would be "allow n subsequent repetitions of an experiment" then the complexity of the algorithm is at least $O(k^n)$, where k is the number of outcomes in the experiment.

4. Solving Linear MDNs

The two methods for solving mdns discussed above restricted the number of plans under consideration. In this section, we introduce another solving method, incremental pruning, that considers all plans, but constrains the value function $V(h)$ of the mdn to a linear one. The strength of this method compared to the former two methods is that not all plans have to be investigated to solve the mdn: an optimal strategy is constructed directly.

4.1. Linear MDNs as POMDPs

Mdns with a linear value function (which will be called linear mdns) are in fact a special case of partially observed Markov decision processes (pomdps). Linear value functions have the form $V(h) = \sum_i U(h)_i + K(h)$. With this form it is possible to distribute the total pay-off value of an episode as partial rewards over the individual steps. At every step but the last one, the immediate reward is the negative cost of the action performed, or of the experiment performed, being $-c_i(a)$ or $-c_i(e)$, respectively. At the last step the reward is the (weighted) expected utility $\sum_i U(h)_i$. The sum of these immediate rewards is the total pay-off value of the episode.

Because linear mdns are pomdps, algorithms for general pomdps are also applicable to linear mdns, although it might not always be such a good idea to do so. Algorithms for general pomdps are designed to deal with difficulties that occur often in normal pomdps but that are less frequent in mdns.

General pomdps (see [1, 6]) are characterized by three discrete and infinite sets and three functions. The sets are: a set of states of the world W , a set of actions A , and a set of observations Z . The functions are: a state-transition probability function $P(w^0|jw; a)$, an observation probability function $P(z|jw; a)$, and a real-valued reward function $r(a; w; w^0; z)$, with $a \in A$; $w; w^0 \in W$ and $z \in Z$ (The reward function can be reduced to a function $r^a(w)$ when rewards do not depend on the observation z or the next state w^0 , as is the case in [1].)

In regard to general pomdps, linear mdns have four special properties: (1) separation of actions and observations, (2) fixed rewards, (3) partition of the state space, and (4) end-state utilities. These properties cause mdns to have a particular structure that can be exploited for efficiency. Below, each property will be explained in detail.

Separation of actions and observations – In standard pomdps every action is immediately followed by an observation. The mdn framework separates the actions from the observations: actions do not lead to observations and experiments do not change the state of the world. A linear mdn is converted into an equivalent pomdp by letting all actions be followed by a fixed single dummy-observation z_0 so that no information is gained by the actions ($P(z_0|jw; a) = 1$), and by using a unit function as state-transition probability for every experiment ($P(w^0 = w|jw; e) = 1$). Furthermore, the observations that can occur at an experiment are restricted to the possible outcomes of that experiment ($P(z|jw; e) = 0$ if $z \notin Z_e$).

The separation of actions and experiments reduces the complexity of computations needed to obtain state probability vectors. In case of an action only a state transition matrix has to be applied; in case of an experiment only Bayes' rule has to be used.

Fixed rewards – The rewards that are obtained at each step in a general pomdp depend on the actual state, the action performed at that step, the observation obtained, and the next state. Often the rewards are a function of a subset of these four factors. The framework for mdns reduces the reward function to a minimum: the rewards at each step, except for the last one, only depend on the action or experiment performed at that step. They are equal to the negative costs of that action or experiment ($r(a; w; w^0; z) = -c(a) = -\sum_i K(a)_i$; and $r(e; w; w^0; z) = -c(e) = -\sum_i K(e)_i$).

Partitioned state space – The most apparent difference between mdns and general pomdps is the partition of the state space. The states in a mdn are partitioned into state groups, such that all probability mass is always concentrated in one group. Even more

important is the fact that the availability of every action and experiment is confined to one group.

End-state utilities – The pomdp “rewards” in mdns are defined as costs of actions and experiments. The real rewards are obtained after the final step when an end state is reached. These rewards are called utilities in mdns and the values of these utilities are determined by the end state. The fact that utilities can only be obtained in end states means that precisely those strategies that always lead to the end group are feasible.

4.2. Incremental pruning

An algorithm that fits naturally to the kind of special properties of mdns described above, is the incremental-pruning algorithm of Zhang and Liu [1, 9]. The basic step in the derivation of this algorithm, according to [1], is to split the value function:

$$V^0(x) = \max_{a \in A} \sum_{w \in W} r^a(w)x(w) + \gamma \sum_{z \in Z} P(z|x; a)V(x_z^a) \quad (4.1)$$

into simpler parts:

$$V^0(x) = \max_{a \in A} V^a(x) \quad (4.2)$$

$$V^a(x) = \sum_{z \in Z} V_z^a(x) \quad (4.3)$$

$$V_z^a(x) = \frac{\sum_{w \in W} r^a(w)x(w)}{|Z|} + \gamma \sum_{z' \in Z} P(z'|z; a)V(x_{z'}^a) \quad (4.4)$$

and to solve them separately:

$$S^0 = \text{purge} \left[\sum_{a \in A} S^a \right] \quad (4.5)$$

$$S^a = \text{purge} \left[\sum_{z \in Z} S_z^a \right] \quad (4.6)$$

$$S_z^a = \text{purge} (f_z^a(\cdot; a; z) \otimes S_g) \quad (4.7)$$

$$f_z^a(\cdot; a; z)(w) = \frac{(1-\gamma) \sum_{w^0 \in W} r^a(w^0) + \gamma \sum_{w^0 \in W} P(w^0|z; a) f^a(w^0)}{P(z|w^0; a)P(w^0|z; a)} \quad (4.8)$$

The solution $V(x)$ of a pomdp is expressed in a set S of $|W|$ -sized vectors \mathbb{R} , each vector representing a separate strategy. In the formulas above, the notation x is used to denote the actual belief state, which is a probability vector defined on the states of the world, where $x(w)$ denotes the probability of the world being in state w . The notation x_z^a is used to denote the revised belief state after performing action a and observing z . The constant γ is a discount rate that is only used in infinite horizon cases, otherwise γ is 1.

The function `purge` is an algorithm that reduces a given set of $|W|$ -sized vectors to those vectors \mathbb{R} for which $\mathbb{R} \cdot x$ is maximal for some belief state x . The “incremental pruning” in a more strict sense is laid down in the solution of equation (4.6):

$$S^a = \text{purge} \left[\sum_{z \in Z} S_z^a \right] \\ = \text{purge} (S_{z_1}^a \otimes S_{z_2}^a \otimes S_{z_3}^a \otimes \dots \otimes S_{z_j}^a) \quad (4.9) \\ = \text{purge} (\dots (\text{purge} (S_{z_1}^a \otimes S_{z_2}^a) \otimes S_{z_3}^a) \dots \otimes S_{z_j}^a) \quad (4.10)$$

The operator \otimes has the following meaning: if S_1 and S_2 are sets of vectors of equal length, then $S_1 \otimes S_2$ is the set of vectors: $\{s_1 + s_2 \mid s_1 \in S_1; s_2 \in S_2\}$.

4.3. Specializing incremental pruning for MDNs

There are in fact three types of decisions in mdns: (1) the decision to perform an action a , (2) the decision to perform an experiment e , and (3) the decision to stop (d_s). The last decision is available whenever the world is in the end group G_u . Each one of the three types of decision has its own solution method: for stop-decisions, just the utility vector has to be used; for action-decisions, a transition matrix has to be applied; and for experiment-decisions, all outcomes have to be inspected. Hence, formula (4.2) can be rewritten into (4.11). Each of the partial solutions V^u , V^a , and V^e can be worked out separately according to the type of decision:

$$V^0(x) = \max_{w \in W} \left[\max_{a \in A(G(x))} V^a(x) \vee \max_{e \in E(G(x))} V^e(x) \right] \quad (4.11)$$

Because the probability mass in a belief state x is always concentrated in one state group, it makes sense to use the expression $G(x)$ to denote that state group. As a result, $A(G(x))$ and $E(G(x))$ form the decision options that are available at belief state x .

The solution for V^u is straightforward: the reward when stopping is the weighted expected utility $\sum_j u_j x_j$:

$$V^u(x) = \sum_j u_j x_j \quad (4.12)$$

At action-decisions, we just have to apply the appropriate transition matrix A_a to compute x_z^a in formula (4.4). The rewards $r^a(s)$ are laid down in the costs $j \in K(a)$.

$$V^a(x) = V(x_z^a) \quad j \in K(a) = V(A_a x) \quad j \in K(a) \quad (4.13)$$

For experiment-decisions, we need the following formula:

$$V^e(x) = \sum_z P(zjx; e) V(x_z^e) \quad j \in K(e) \quad (4.14)$$

The posterior vector x_z^e can be computed by application of the sensitivity matrix E_e :

$$x_z^e = \frac{d(E_e^z):x}{P(zjx; e)} \quad (4.15)$$

We use $d(E_e^z)$ to denote the diagonal matrix whose trace is equal to row z in matrix E_e . This notation seems awkward, but because $V(x)$ is a linear function, the term $P(zjx; e)$ will cancel out and only $d(E_e^z)$ remains, so $V^e(x)$ becomes:

$$V^e(x) = \sum_z V(d(E_e^z):x) \quad j \in K(e) \quad (4.16)$$

These formulas result in the following incremental pruning (ip) scheme:

$$S^0(G) = \text{purge}^i S^d(G) \quad (4.17)$$

$$S^d(G) = \sum_{j \in K(G)} S^u(G) \cup S^a(G) \cup S^e(G) \quad (4.18)$$

$$S^u(G) = \begin{cases} f_{j \in K(G)} u_j & \text{if } G = G_u \\ \emptyset & \text{if } G \neq G_u \end{cases} \quad (4.19)$$

$$S^a(G) = \sum_{a \in A(G)} S^a(G) \quad (4.20)$$

$$S^e(G) = \sum_{e \in E(G)} S^e(G) \quad (4.21)$$

$$S^a(G) = \text{purge}(f_{j \in K(G)} A_{aj} \sum_{j \in K(a)} S^{p(a)}(G^0) g_j) \quad j \in K(a) \quad (4.22)$$

$$S^e(G) = \text{purge}(\sum_{z \in Z_e} S_z^e(G)) \quad j \in K(e) \quad (4.23)$$

$$S_z^e(G) = \text{purge}(f_{j \in K(G)} d(E_e^z) \sum_{j \in K(e)} S^{p(e)}(G) g_j) \quad (4.24)$$

(The notation $S^{p(\cdot)}$ will be explained in the next section.) As is apparent from these formulas, a separate ip scheme exists for every state group G . In each step during the value-iteration process, this ip scheme can be performed for every group independently from the schemes of other groups. The complexity of an ip scheme is determined by the size of its state group, the number of actions and experiments attached to this group, and the number of outcomes that each experiment has.

4.4. Linking IP-scheme results

The ip scheme for a state group is linked to the results of ip schemes in previous value-iteration steps. In formula (4.22) a link is made to the ip scheme of the outgoing group $G^{\text{out}}(a)$ of action a through the expression $S^{p(a)}(G^{\text{out}}(a))$. A similar link is made through $S^{p(e)}(G)$ in formula (4.24).

The ip scheme can only be performed for a state group if that state group can have probability mass in the actual value-iteration step. At the first value-iteration step, for instance, only the end group G_u can have probability mass, because every strategy in an mdn should finish in the end group for all episodes. At the second value-iteration step only those state groups can have positive probability mass when the end group is reachable in one step through an action in the mdn graph. In general, at the n -th value-iteration step, only those groups can have positive probability mass, and therefore have their ip scheme performed, if the end group is reachable in exactly $n-1$ steps in the mdn graph. If a link has to be made from one ip scheme to another in a previous step, the algorithm has to look for the most recent step in which the target ip scheme is available. In $S^{p(a)}(G^{\text{out}}(a))$, the function $p(a)$ denotes the "previous" step, which is the most recent value-iteration step in which an ip scheme exists for the group $G^{\text{out}}(a)$.

```

MdnIncPrune (M = hA; E; Z; W; Gi, maxHor)
1.  S0(;)  $\bar{A}$  ;, S0(Gu)  $\bar{A}$  {(di, s:u)}
2.  for h  $\bar{A}$  1 to maxHor:
3.    Sh(;)  $\bar{A}$  ;
4.    for all G  $\in$  G:
5.      p  $\bar{A}$  max(q j q < h; Sq(G)  $\notin$  ;))
6.      if p  $\notin$  ? then
7.        Sh(G)  $\bar{A}$  purge(Sh(G) [  $\bar{A}$ (Sp(G); d0)])
8.        for all e  $\in$  E(G):
9.          S0  $\bar{A}$  purge(f@:d(Ee1)j@ 2Sp(G)g)
10.         for z  $\bar{A}$  2 to jZej:
11.           S0  $\bar{A}$  purge(S0  $\odot$ 
12.             purge(f@:d(Eez)j@ 2Sp(G)g))
13.           Sh(G)  $\bar{A}$  purge(Sh(G) [  $\bar{A}$ (S0; de)])
14.         for all a  $\in$  A(G):
15.           G0  $\bar{A}$  Gout(a)
16.           p  $\bar{A}$  max(q j q < h; Sq(G0)  $\notin$  ;))
17.           if p  $\notin$  ? then
18.             S0  $\bar{A}$  purge(f@:Aaj@ 2Sp(G0)g)
19.             Sh(G)  $\bar{A}$  purge(Sh(G) [  $\bar{A}$ (S0; da)])
20.         If Sh(G0)  $\notin$  ; and d = d0
21.         for all decisions in Sh(G0) then Stop.

```

Figure 4.1: Implementation of the incremental-pruning algorithm attempting to solve the mdn M in at most maxHor steps.

4.5. The incremental-pruning algorithm for MDNs

Figure 4.1 gives the pseudocode for an implementation of the incremental-pruning algorithm for mdns. The value-iteration steps are indicated here by the horizon variable h . A separate ip scheme is performed at every step for each state group. The results of the ip schemes for state group G at step h are represented by the set $S^h(G)$ of tuples $hd, @i$.

At line 7 a new, auxiliary decision is introduced: the wait-decision d_0 . This decision has no meaning in the mdn itself, but is used in this algorithm to implement a stopping mechanism. The wait-decision is a decision to keep everything exactly the same for one step. The set $S^h(G)$ is filled at line 7 with the contents of the previous vector set for this group. The algorithm stops at line 19 when no decision can give a higher expected value than the wait-decision for any belief state in the start group G_0 , which means that $S^h(G_0)$ only contains wait-decisions. The wait-decision has another effect: it is now only necessary to add a stop-

decision to $S^0(G_u)$. The other stop-decisions that are mentioned in formula (4.19) of the ip scheme are replaced by wait-decisions added to $S^h(G_u)$. Function $\bar{A}(S; d)$ at lines 7, 12 and 18 returns a set that contains all entries of S having the decision set to d and having $@$ diminished with the costs of the decision d .

The purge routine called at lines 9, 11, 12, 17, and 18 is implemented by a slightly altered version of Lark's Filter algorithm as described in [1, 10]. The iterative purging process as described at line 11 is in fact the most simple incremental-pruning method possible and could well be replaced by one of the more efficient methods (see [1]).

5. Time complexity of the incremental-pruning algorithm for MDNs

We are curious to see whether a decrease of the time complexity is caused by the special structure of mdns. More specifically, we would like to investigate the role of the division of the state space into state groups. Hence, we are not aiming at obtaining a general formulation of the time complexity of the incremental-pruning algorithm for mdns,

Figure 5.1 provides the time complexity for each line of the algorithm given in Figure 4.1. The time complexities are based on the following assumptions:

- ² The number of value-iteration steps needed is at most \hat{h} .
- ² The size of all state groups is at most \hat{s} .
- ² The number of actions per group is at most \hat{a} .
- ² The number of experiments per group is at most \hat{e} .
- ² The number of outcomes per experiment is at most \hat{r} .
- ² The number of \otimes -vectors produced per group per value-iteration step is at most \otimes .
- ² The time complexity of the purge algorithm is at most $T(\text{purge}(\otimes; \hat{s}))$.

From Figure 5.1 we can derive the following formula for the overall complexity:

$$T \leq 1 + \hat{h} \cdot [1 + \sum_{j \in G} [\hat{h} + 1 + \otimes \cdot (\hat{s} + 1) + T_p + \hat{e} \cdot f(\otimes; \hat{s} + T_p) + (\hat{r} - 1) \cdot (\otimes \cdot \hat{s} + \otimes^2 \cdot \hat{s} + 2 \cdot T_p) + \otimes \cdot (\hat{s} + 1) + T_{pg} + \hat{a} \cdot f_2 + \hat{h} + \otimes \cdot \hat{s}^2 + T_p + \otimes \cdot (\hat{s} + 1) + T_{pg}] + \otimes] \quad (5.1)$$

(The symbol 'Tp' is an abbreviation of $T(\text{purge}(\otimes; \hat{s}))$.) This formula can be rearranged by collecting all terms in \hat{h} , \otimes , \hat{s} , and T_p :

1.	1
2.	$\hat{h} \in \{$
3.	1
4.	$\sum_{j \in G} \{$
5.	\hat{h}
6.	1
7.	$\otimes \cdot (\hat{s} + 1) + T(\text{purge}(\otimes; \hat{s}))$
8.	$\hat{e} \in \{$
9.	$\otimes \cdot \hat{s} + T(\text{purge}(\otimes; \hat{s}))$
10.	$(\hat{r} - 1) \in \{$
11.	$\otimes \cdot \hat{s} + \otimes^2 \cdot \hat{s} + 2 \cdot T(\text{purge}(\otimes; \hat{s})) \}$
12.	$\otimes \cdot (\hat{s} + 1) + T(\text{purge}(\otimes; \hat{s})) \}$
13.	$\hat{a} \in \{$
14.	1
15.	\hat{h}
16.	1
17.	$\otimes \cdot \hat{s}^2 + T(\text{purge}(\otimes; \hat{s}))$
18.	$\otimes \cdot (\hat{s} + 1) + T(\text{purge}(\otimes; \hat{s})) \}$
19.	$\otimes \}$

Figure 5.1: Complexity per line of the algorithm given in Figure 4.1.

$$T \leq 1 + \hat{h} + \hat{h} \cdot \otimes + (\hat{h} \cdot \sum_{j \in G}) \cdot [1 + 2 \cdot \hat{a}] + [1 + \hat{a}] \cdot \hat{h} + [1 + \hat{e} + \hat{a}] \cdot \otimes + [\hat{e} \cdot (\hat{r} - 1) + \hat{a} + 1] \cdot \otimes \cdot \hat{s} + [\hat{e} \cdot (\hat{r} - 1)] \cdot \otimes^2 \cdot \hat{s} + [\hat{a}] \cdot \otimes \cdot \hat{s}^2 + [2 \cdot \hat{e} \cdot \hat{r} + 2 \cdot \hat{a} + 1] \cdot T_p \quad (5.2)$$

The square-bracketed terms in this formula are all bounded by:

$$\hat{d} = 2 \cdot \hat{e} \cdot \hat{r} + 2 \cdot \hat{a} + 1 \quad (5.3)$$

So, the formula can be rewritten into:

$$T \leq 1 + \hat{h} + \hat{h} \cdot \otimes + (\hat{d} \cdot \hat{h} \cdot \sum_{j \in G}) \cdot [1 + \hat{h} + \otimes + \otimes \cdot \hat{s} + \otimes^2 \cdot \hat{s} + \otimes \cdot \hat{s}^2 + T_p] \quad (5.4)$$

When the state groups in an mdn are of comparable size and each group has a comparable number of actions, experiments, and outcomes per experiment, then the following approximations can be made:

$$\hat{d} = \frac{D}{jGj} \quad (5.5)$$

$$\hat{s} = \frac{jWj}{jGj} \quad (5.6)$$

where D is defined by:

$$D = 2:jEj: \max_{e \in E} jR_ej + 2:jAj + 1 \quad (5.7)$$

This leads to the following formula for the time complexity:

$$T < 1 + \hat{h} + \hat{h} \cdot \textcircled{*} + (D:\hat{h}):f1 + \hat{h} + \textcircled{*} + \textcircled{*} \cdot \frac{jWj}{jGj} + \textcircled{*}^2 \cdot \frac{jWj}{jGj} + \textcircled{*} \cdot \frac{jWj^2}{jGj^2} + T_p g \quad (5.8)$$

which ...nally can be rewritten into:

$$T \leq [1 + \hat{h} + \hat{h} \cdot \textcircled{*} + D:\hat{h}:(1 + \hat{h} + \textcircled{*})] + \frac{[D:\hat{h}:(\textcircled{*} + \textcircled{*}^2):jWj]}{jGj} + \frac{[D:\hat{h}:\textcircled{*}:jWj^2]}{jGj^2} + [D:\hat{h}]:T_p \quad (5.9)$$

On base of this formula, the time complexity of the incremental-pruning algorithm for mdns can be expressed as a function of the number of state groups, jGj:

$$T(jGj) = O(K_1 + \frac{K_2}{jGj} + \frac{K_3}{jGj^2} + K_4:T(\text{purge}(\textcircled{*}, \frac{jWj}{jGj}))) \quad (5.10)$$

where $K_1::K_4$ are the terms that do not depend on the number of groups. The time complexity of the purge routine is a function of the number of vectors ($\textcircled{*}$) and the size of the vectors ($jWj=jGj$). This function depends highly on the implementation of the Filter algorithm, but for the size of the vectors the time complexity will be more than linear.

To conclude, an increase in the number of groups jGj will either have no effect on a term in expression (5.10) or will lower the value of a term. The total time complexity will therefore decrease if the number of state groups increases. Going from a general pomdp with only one state group to an mdn with multiple state groups will therefore decrease the total time complexity of the incremental-pruning algorithm.

6. Example

We will use a small example to show the working of mdns. Assume an agent has to go through 9 different subsequent stadia before reaching a final rewarding stadium. In every stadium the agent can be in one of two states (a good state and a bad one) but the actual state is unknown to the agent.

Going from one stadium to the next one, the agent can try to improve his situation by investing 25 units, having a probability of 30% to reach a good state from a bad state, but also having a probability of 20% to arrive in a bad state from a good one. The agent can decide to do nothing instead of investing, in which case he has only a probability of 10% to change his state from a bad one to a good one and also a probability of 10% to change from a good state into a bad one.

Investment actions are named a, b, c, \dots and the do-nothing actions a^0, b^0, c^0, \dots . At every stadium, except the last one, the agent can use an experiment (named A, B, C, ...) with two outcomes to gain information on the actual state. The experiments cost 10 units each but are erroneous in 5% of the cases. The utility for reaching the bad state in the last stadium is 0, and is 1 for reaching the good state. The value function used is: $V(h) = \alpha U(h) + K(h)$, which is linear. The structure of the mdn that models this problem is given in figure 6.1.

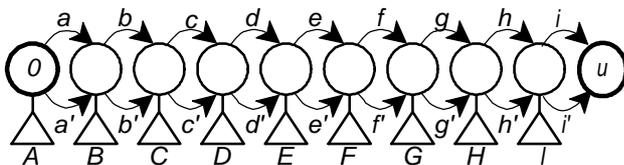


Figure 6.1: Graphical representation of the 10-stadia example.

Depending on the weight α , the incremental-pruning algorithm finds the following optimal strategies:

- $s_1 : ha^0b^0c^0d^0e^0f^0g^0h^0i^0 \} i$
- $s_2 : ha^0b^0c^0d^0e^0f^0g^0h^0I \} i^0 \} i$
- $s_3 : ha^0b^0c^0d^0e^0f^0g^0HhI \} i^0 \} h^0i^0 \} i$
- $s_4 : ha^0b^0c^0d^0e^0f^0GgHhI \} i^0 \} h^0i^0 \} g^0h^0I \} i^0 \} i$
- $s_5 : ha^0b^0c^0d^0e^0f^0GgHhI \} i^0 \} h^0I \} i^0 \} i^0 \} g^0h^0I \} i^0 \} i$

In s_2 we see, for example, the application of experiment I. It decides on a continuation with i or i^0 . The solution sets of strategies that are optimal for some region of the prior belief space are for different levels of α :

- $\alpha = 1 : fs_{1g}$ $\alpha = 300 : fs_{2;S_3g}$
- $\alpha = 200 : fs_{1g}$ $\alpha = 325 : fs_{3g}$
- $\alpha = 250 : fs_{1;S_2g}$ $\alpha = 475 : fs_{3;S_4g}$
- $\alpha = 275 : fs_{2g}$ $\alpha = 500 : fs_{3;S_4;S_5g}$

So, the higher the weight factor on the utilities, the more the agent is willing to invest in experiments. The understanding that investing is a good (even optimal) strategy, starts in the last stadium and progresses forwards in relation to an increase of the weight α . The computations took at most 7 seconds on a Pentium II, 300 Mhz personal computer.

7. Conclusion

The new technique for modeling decision planning in uncertain environments, *mdn*, adequately combines the flexibility of *pomdps* with the semantics of influence diagrams. Having shown that linear *mdns* are *pomdps*, we adapted the most efficient algorithm for solving *pomdps* known today i.e., the incremental pruning algorithm, to make it suitable for linear *mdns*. We showed that the incremental-pruning algorithm adapted for *mdns* is more efficient in terms of time complexity than that for general *pomdps*.

Acknowledgement

The research is supported by the Dutch organization for air (Dutch: Lucht) and Water environment Information systems (Iwi).

8. References

- [1] Cassandra, A.R., Littman, M.L., and Zhang, N.L. (1997). Incremental Pruning: A Simple, Fast, Exact Method for Partially Observable Markov Decision Processes. Proceedings of the thirteenth conference on Uncertainty in Artificial Intelligence, pp. 54-61. Morgan Kaufmann Publ., San Mateo, CA. ISBN 1-55860-485-5.
- [2] Donkers, H.H.L.M. and Hasman, A. (1988). Strategy, een Nieuwe Kijk op Beslisbomen. Proceedings of the 8st Dutch Medical Informatics Congress, pp. 65-69. VMBI, Utrecht, The Netherlands. ISBN 90-5005-022-0.
- [3] Donkers, H.H.L.M. (1997). Markov Decision Networks. Technical Report CS 97-04. Universiteit Maastricht, Maastricht, The Netherlands.
- [4] Donkers, H.H.L.M, Uiterwijk, J.W.H.M., and Herik, H.J. van den (1998). Markov Decision Networks. Submitted to ECAI'98.
- [5] Donkers, H.H.L.M, Uiterwijk, J.W.H.M., and Herik, H.J. van den (1998). Solving Markov Decision Networks using Incremental Pruning. Submitted to UAI'98.
- [6] Lovejoy, W.S. (1991). A Survey of Algorithmic Methods for Partially Observed Markov Decision Processes. Annals of Operations Research, 28(1-4), pp. 47-66. ISSN 0254-5330.
- [7] Olmsted, S.M. (1983). On Representing and Solving Decision Problems. Ph.D. Dissertation Engineering-Economic Systems. Stanford University, Stanford, CA.
- [8] Pomerol, J.C. (1997). Artificial Intelligence and Human Decision Making. European Journal of Operations Research, 99, pp. 3-25. ISSN 0377-2217.
- [9] Zhang, N.L. and Liu, W (1996). Planning in Stochastic Domains: Problem Characteristics and Approximation. Technical Report HKUST-CS96031, Hong Kong University of Science and Technology.
- [10] Zhang, N.L. and Liu, W (1997). A Model Approximation for Planning in Partially Observable Stochastic Domains. Journal of Artificial Intelligence Research, 7, pp. 199-230.