

The Complexity of Model Aggregation

Judy Goldsmith

Dept. of Computer Science
University of Kentucky
763 Anderson Hall
Lexington, KY 40506-0046
goldsmi@cs.uky.edu

Robert H. Sloan

Dept. of EE & Computer Science
University of Illinois at Chicago
851 S. Morgan St. Rm 1120
Chicago, IL 60607-7053
sloan@cecs.uic.edu

Abstract

We show that the problem of transforming a structured Markov decision process (MDP) into a *Bounded Interval MDP* is coNP^{PP} -hard. In particular, the test for ϵ -homogeneity, a necessary part of verifying any proposed partition, is coNP^{PP} -complete. This indicates that, without further assumptions on the sorts of partitioning allowed or the structure of the original propositional MDP, this is not likely to be a practical approach. We also analyze the complexity of finding the minimal-size partition, and of the k -block partition existence problem. Finally, we show that the test for homogeneity of an exact partition is complete for $\text{coNP}^{\text{C}^{\text{P}}}$, which is the same class as coNP^{PP} .

All of this analysis applies equally well to the process of partitioning the state space via Structured Value Iteration.

Introduction

Markov decision processes (MDPs, formally defined below) are ubiquitous in AI and in the world of mathematical modeling. Related research concerns learning models and/or policies for complex systems, and developing algorithms and heuristics for planning and intelligent control. Standard AI applications include planning, robot control, medical systems, logistics, and many others. MDPs model controllable stochastic processes: there is a set of *states*; a controller chooses among some number of *actions*, each of which has an associated probability matrix of state transitions; associated with each state and action pair is a *reward*. The basic goal, given such a model, is to find a strategy or *policy* for choosing actions that maximizes the total expected reward over some specified time horizon.

There are two reasons why finding good policies for MDPs is computationally complex. The first is the underlying complexity of finding optimal policies. The second is the sheer size of most interesting models of controlled stochastic systems, especially AI planning problems, where one often considers a “size n ” problem to consist of an $\Omega(2^n)$ state space described by n *fluents*

or variables, and some appropriate data structure representing the dependence relations of each each fluent to the states of some small number of fluents at the current or previous stage. One can view this sort of model as an attempt to reduce the sizes of the representations of MDPs. Unfortunately, simply reducing the size of the representation is not good enough. Although there are these factored (also called “structured” or “propositional”) representations (Bayes’ nets and others) of size n that represent 2^n -state MDPs, the problem of deciding whether a good plan exists is correspondingly exponentially harder than for the unfactored representations. (For details of this blow-up, see (Goldsmith & Mundhenk 1998; Mundhenk *et al.* 1999).) These factored representations are used because they often model human understanding of the system (as derived from a domain expert), even if the compression of transition tables to simpler or at least smaller forms does not automatically buy any computational speed-up.

One approach to further reducing the representation of a succinctly represented MDP is to group states together into meta-states, where all elements of a meta-state behave exactly or approximately the same with respect to the reward function and all actions. These can be described as aggregate approximate models. We use the terminology of Givan *et al.* (Dean, Givan, & Leach 1997; Dean & Givan 1997; Givan, Leach, & Dean 1997; Givan & Dean 1997) of *Bounded Interval MDPs* (BMDPs). These are MDP-like models where the transition probabilities and rewards are replaced by intervals. Work on planning algorithms for such systems by Givan *et al.* is reported in (Givan, Leach, & Dean 1997) and by Harmanec in (Harmanec 1999). Givan *et al.* give an (exponential worst case) algorithm for partitioning a traditional MDP state space into aggregate states to form a BMDP.

Givan and Dean have shown that, for factored MDPs, the problem of finding an optimal, that is, coarsest such partition for a given interval diameter ϵ is NP-hard (Givan & Dean 1999). We show here that it is in fact much harder (assuming that $\text{P} \neq \text{NP}^{\text{PP}}$). It is not surprising that finding good exact or approximate partitions is difficult. Finding an optimal policy for a finite-horizon succinctly represented MDP is PSPACE-hard (Mund-

ary policies for succinctly represented MDPs is EXP-hard (Lusena, Goldsmith, & Mundhenk 1999). This means that if there is an algorithm that, on input of a succinctly represented MDP, outputs a stationary policy with value v , such that v minus the optimal stationary policy value for that MDP is less than ϵ , then that algorithm can be used to solve EXP-hard problems. In particular, this guarantees that such an algorithm cannot run in time polynomial in the size of the input. Therefore, we cannot expect to find a polynomial-time algorithm that yields significantly smaller representations of succinctly represented MDPs. On the other hand, the results on asymptotic complexity do not rule out particular instances or classes of instances for which a fast and guaranteed good approximation algorithm exists, nor does it rule out fast instances for a general algorithm.

A related approach to Givan *et al.*'s, is given by Boutilier, Dearden, *et al.* in (Boutilier, Dearden, & Goldszmidt 1995; Boutilier & Dearden 1996; Boutilier, Dearden, & Goldszmidt 1999), using structured policy representations. The major contribution of that work is to adapt Value Iteration (Bellman 1957) or Modified Policy Iteration (MPI) (Puterman & Shin 1978) to work on Dynamic Bayes net (2TBNs, defined below) policies represented as trees. A policy tree consists of a set of Boolean formulas that partition the state space, and for each element of the partition, an action. This means that when the specification of a state satisfies a particular formula, the corresponding action is applied.

We argue that the extension step in their algorithms is exactly as hard as the stability tests for BMDPs.

In fact, we argue that any partition of the state space of a factored or structured MDP into states that behave "approximately the same" will lead to an NP^{NP}-hard stability test.

Definitions

A MDP describes a controlled stochastic system by its states and the consequences of actions on the system. It is denoted as a tuple $M = (S, s_0, \mathcal{A}, t, r)$, where

- S and \mathcal{A} are finite sets of *states* and *actions*,
- $s_0 \in S$ is the *initial state*,
- $T : S \times \mathcal{A} \times S \rightarrow [0, 1]$ is the *state transition function*, where $t(s, a, s')$ is the probability to reach state s' from state s on action a (where $\sum_{s' \in S} t(s, a, s') \in \{0, 1\}$ for $s \in S, a \in \mathcal{A}$),
- $r : S \times \mathcal{A} \rightarrow \mathbf{Z}$ is the *reward function*, where $r(s, a)$ is the reward gained by taking action a in state s .

2TBNs

There are many factored representations for MDPs available. We use 2-Phase Temporal Bayes' Networks (2TBNs) as a convenient model for our examples.

A 2TBN consists of a set of variables, or fluents, a set of actions, and a reward function. The effect of

each action on each fluent at time $t + 1$ is determined by the states of a (usually small) subset of the fluents at time t and perhaps $t + 1$. These dependencies are modeled as a directed acyclic graph, where a directed edge indicates dependence. (Although one could build cyclic such graphs, the definitions rule out that case.) The nodes of the graph consist of two sets. One set represents the set of fluents at time t , and the other, the fluents at time $t + 1$. Edges can go from the first set to the second (asynchronous), or within the second set (synchronous).

In addition to this graph, for each time $t + 1$ node (intuitively, for each fluent) and each action, there is a data structure (formula, table, or tree, usually) that represents the effects of that action on that fluent, as a function of that node's parent nodes.

If we assume that each fluent is binary, then a 2TBN with n fluents models a MDP with 2^n states. (If a fluent can have more settings, there can be even more states.) When the dependencies specified by the 2TBN are small enough, the size of the representation may be polynomial in n instead of 2^n . These are the 2TBNs of interest; most examples of 2TBNs in the planning literature have such sparse dependencies. Mundhenk *et al.* showed that the representational savings do not lead to any computational savings, because there is an exponential jump in worst-case complexity as you go from uncompressed to compressed representations. However, models such as 2TBNs are certainly worth using, both for their readability and for computational reasons. For instance, if the effects of actions are represented as arithmetic decision diagrams, this allows use of commercial optimizing software and brings a corresponding increase in computational ability (Hoey *et al.* 1999). (For a fuller discussion of Bayesian networks in planning, see (Boutilier, Dean, & Hanks 1999; Blythe 1999).)

Model Reduction

We next present the approach of Givan *et al.* (Dean, Givan, & Leach 1997; Givan, Leach, & Dean 1997; Givan & Dean 1997) for constructing bounded interval MDPs (BMDPs) out of MDPs. This holds great promise for reducing the time needed to find an approximately optimal policy for an MDP, for those MDPs that yield significantly smaller BMDPs. Unfortunately, the "model reduction techniques" that Givan *et al.* present are potentially quite slow.

A BMDP differs from a standard MDP in that rewards are expressed as a range of possible values, and each transition probability is an interval within $[0, 1]$. A BMDP can be interpreted as a family of MDPs such that each MDP in the family has the same state set and the same action set as the BMDP, and for each MDP, each reward and transition falls within the range given for the BMDP. Alternatively, each state of the BMDP may represent an aggregate state of some base MDP. The advantage of a BMDP over an MDP arises when the number of states of the BMDP is significantly

The basic approach given in (Dean & Givan 1997) for constructing a BMDP from an MDP is to first partition the states of the MDP by reward, and then to iterate the following: while there are aggregate states A and B such that there is an action a and states s and s' in A so that

$$\left| \sum_{t \in B} T(s, a, t) - \sum_{t \in B} T(s', a, t) \right| > \epsilon \quad (1)$$

then split A into $A_1 \dots A_k$ so that for each A_i and all s and s' in A_i ,

$$\left| \sum_{t \in B} T(s, a, t) - \sum_{t \in B} T(s', a, t) \right| \leq \epsilon.$$

We say that the initial aggregate state A is ϵ -unstable with respect to initial aggregate state B if inequality 1 holds for A and B : the opposite condition is called ϵ -stability. The goal is to produce a partition that is ϵ -homogeneous, that is, ϵ -stable with respect to all pairs A and B of aggregate states.

Givan and Dean show that the BMDP induced by an ϵ -partition is a "close enough" approximation to the original MDP, namely that an appropriate version of value iteration gives a policy for the BMDP that defines a reasonable approximation to the optimal policy for the original MDP. (The closeness of the approximation depends, of course, on ϵ .) Unfortunately, they also show that the test for whether there is a splittable aggregate state A is NP-hard in the case of a factored (propositional) state space (Givan & Dean 1999). Thus, there is no obvious polynomial-time stopping criterion for this algorithm in the factored case. Our work presents an even more precise picture of this complexity.

On the one hand, we have exactly pinpointed the complexity of the problem, which offers the opportunity to apply heuristics designed for exactly such problems. On the other hand, the complexity classes discussed are conjectured to properly contain the class NP (as discussed below). Thus, we have shown that the problem is even worse than NP-hard.

Representing Partitions We make some assumptions here about how partitions will be represented. These are minimal assumptions necessary for there to be any hope that factored representations can help. If they are violated, then the MDP inherently has super-polynomially many "very different" states that must be treated separately. In practice, *stronger* assumptions will typically be made (Givan 1999). Specifically, we assume that every partition is given implicitly (rather than as an enumeration of states into blocks) so that the partition may be represented in size polynomial in the size of the succinctly represented MDP. For instance, if the MDP is represented as a 2TBN, the blocks of the partition could be represented as Boolean formulas on the fluents. Furthermore, in order for there to be an appreciable savings in complexity, we hope that the

number of blocks in the partition is quite small relative to the size of the state space.

If the representation of a partition has size $\Omega(2^n)$, for instance, where n is the number of fluents in the 2TBN, then the algorithm for checking stability is linear in the size of the input, namely the 2TBN plus the partition. This does not, however, represent a win in complexity terms.

Relevant Complexity Classes

For definitions of standard complexity classes, reductions, and results from complexity theory we refer to (Papadimitriou 1994). We use the notation FP to refer to functions computable in time polynomial in the size of the input.

Let M be a nondeterministic polynomial-time Turing machine, and let $Acc_M(x)$ be the number of accepting computations of M on input x , and $Rej_M(x)$ the number of rejecting computations. The class NP is the class of all sets $E_M = \{x : Acc_M(x) \geq 1\}$. One can characterize NP in terms of an existential quantifier: there exists an accepting computation of $M(x)$. The SAT is the canonical complete set for NP: SAT is in NP, and for every set S in NP, there is a polynomial-time reduction from S to SAT. We define $NP = \Sigma_1^P$ and $coNP = \Pi_1^P$. We can then build a hierarchy of classes Σ_k^P and their complements, $\Pi_k^P = co\Sigma_k^P$, by alternations of quantifiers (\exists and \forall). Equivalently, one can use a complete set, A_k , for Σ_k^P as an oracle to define $L \in \Sigma_{k-1}^P$ if and only if there is a nondeterministic polynomial-time oracle Turing machine $M^{(l)}$ such that $x \in L \Leftrightarrow M^{A_k}(x)$ accepts. The *Polynomial Hierarchy* (PH) consists of all classes Σ_k^P and Π_k^P . All these classes can be accepted by polynomial space bounded Turing machines (and thus by exponential time bounded Turing machines), so are contained in PSPACE.

The class PP is the class of all sets $S_M = \{x : Acc_M(x) > Rej_M(x)\}$. Intuitively, S_M is the set of inputs to M that are more likely to be accepted than rejected on any random choice of computation path. We refer to PP computations as *probabilistic polynomial-time computations*, or just probabilistic computations. This does not imply, however, that one can be very sure of the answer one gets on a single or several computations. If one wants to increase confidence by repeating trials, one must put stronger strictures on the probabilities, for instance that more than 3/4 of the computations are correct. That defines a different complexity class, BPP, which is apparently less powerful than PP. All PP functions are computable in polynomial space. Furthermore, Toda showed that computing PP was at least as hard as computing any set in the polynomial hierarchy (Toda 1991). That is, $PH \subseteq P^{PP} \subseteq PSPACE$.

The function class #P is the set of functions f such that for some nondeterministic polynomial-time Turing machine M , $f(x) = Acc_M(x)$. One can see that knowing the precise number of accepting computations would answer NP questions and PP questions. It is less obvious, but true, that being able to look up the

precise number of accepting computations would answer any membership question for the polynomial hierarchy (Toda 1991).

One can express the Polynomial Hierarchy in terms of \exists and \forall quantifiers. One can define a similar hierarchy in terms of \exists , \forall , and *counting quantifiers*, where for any $f \in \text{FP}$, the counting quantifier C_f is defined as $C_f(y)R(x, y) \Leftrightarrow |\{y : |y| \leq p_f(|x|) \wedge R(x, y)\}| \geq f(|x|)$, where p_f is some polynomial. The function f may be implicit; the standard version of the quantifier is just C . The hierarchy built up from the class P using these three quantifiers is called the *Counting Hierarchy*, (CH) and was first defined by Wagner (Wagner 1986). The class CP is exactly the class PP , and NP^{PP} is $\exists\text{CP}$.

The class NP^{PP} consists of sets for which one can guess a proof of membership and *probabilistically check* its correctness. The class NP^{PP} captures the complexity of asking “Is there a good policy for this factored unobservable MDP?” (In other words, “Is there a good linear plan?”) This was first shown in (Goldsmith, Littman, & Mundhenk 1997). While a fast deterministic algorithm for an NP^{PP} -complete problem is almost certainly not going to exist, heuristic algorithms have begun to be developed for such problems (Majercik & Littman 1998b; 1998a; Littman 1999; Littman, Majercik, & Pitassi 2000).

To summarize,

$$P \subseteq \text{NP} \subseteq \text{PH} \subseteq \text{P}^{\text{PP}} \subseteq \text{NP}^{\text{PP}} \subseteq \text{CH} \subseteq \text{PSPACE}.$$

The class $\text{C}_{=}\text{P}$ is the class of languages L such that there exists functions $f \in \#\text{P}$ and $t \in \text{FP}$ such that for all x , $x \in L \Leftrightarrow f(x) = t(x)$ (Wagner 1986). (In other words, there is some polynomial-time nondeterministic TM M such that $M(x)$ accepts on *exactly* $t(x)$ computations.) One can view $\text{C}_{=}$ as an operator, similar to the C quantifier, and define classes in an extended counting hierarchy using this operator. The proper Counting Hierarchy is contained in this extended hierarchy, which is in turn contained in PSPACE . While it is not apparently relevant to this work, it is interesting to note that the class $\text{coC}_{=}\text{P}$ is equal to the nondeterministic quantum complexity class NQP_{C} (Feiner *et al.*; 1999).

To the best of our knowledge, the class $\text{C}_{=}\text{P}$ has not come up in the complexity analysis of planning problems up to now.

Hardness of Approximate Stability Testing

As mentioned before, the approach given in (Dean & Givan 1997) for constructing a BMDP from an MDP is to first partition the states of the MDP by reward, and then to iterate the following: while there are aggregate states A and B such that there is an action a and states s and s' in A so that

$$\left| \sum_{t \in B} T(s, a, t) - \sum_{t \in B} T(s', a, t) \right| > \varepsilon$$

then split A into $A_1 \dots A_k$ so that for each A_i and all s and s' in A_i ,

$$\left| \sum_{t \in B} T(s, a, t) - \sum_{t \in B} T(s', a, t) \right| \leq \varepsilon.$$

Let us formally define the problem FACTORED MDP ε -HOMOGENEITY as follows:

Instance: A factored MDP and a partition of its state space, and a rational number $0 \leq \varepsilon \leq 1$.

Question: Is the partition ε -homogeneous?

Theorem 1 *The FACTORED MDP ε -HOMOGENEITY problem is coNP^{PP} -hard.*

Proof. We give a reduction from EMAJ3SAT , a known NP^{PP} -complete problem (Littman, Goldsmith, & Mundhenk 1998). An instance of EMAJ3SAT consists of a formula ϕ in 3CNF with m clauses over n Boolean variables, x_1, \dots, x_n and an additional parameter k such that $1 \leq k \leq n$. The question is whether there is some assignment to x_1, \dots, x_k such that for that assignment, a majority of assignments to x_{k+1}, \dots, x_n satisfy ϕ .

Given an instance of EMAJ3SAT , we create the following factored MDP M . There are $m + n + 2$ fluents, $c_1, \dots, c_m, x_1, \dots, x_n$, and v_0, v_1 . The description of M will be a little easier to follow if we think of v_0 and v_1 being replaced by a single variable v that takes on the values 0, 1, 2, and 3.

There is a single action, a . If v is 0 or 1, then a maintains the current values of x_1, \dots, x_k , and sets x_{k+1}, \dots, x_n randomly. If v is 2 or 3, then a sets all n of the x_i s at random. Each c_j corresponds to one of the clauses, and a sets each c_j to be true if and only if the *new* values of the three x_i s corresponding to c_j 's literals would cause the corresponding clause to be true. Finally, a changes v as follows: 0 becomes 2, 1 becomes 3, and 2 and 3 both become 0 or 1, each with probability $1/2$.

The reward function assigns reward 100 to any state with $v = 2$ and all $c_j = 1$. Any state with $v = 0$ or $v = 1$ has reward 50; all other states have reward 0.

Finally, we set $\varepsilon = 1/2$.

For any $\varepsilon \leq 1$ (and in fact any $\varepsilon < 25$), the reward function induces an initial partition into three blocks. Let block A be the block where $v = 0$ or 1 (reward 50), block B be the block where $v = 2$ and all $c_j = 1$ (reward 100), and block C be all other states (reward 0).

For our MDP, blocks B and C in the initial partition are both 0-stable with respect to each other block. Now consider the effect of action a on state s in A . If state s has $v = 1$, then a will change v to 3, and send s to a state in C with probability 1. Therefore $\sum_{z \in B} T(s, a, z) = 0$. However, suppose that s has $v = 0$ and for that setting of x_1, \dots, x_k , more than half of the settings of x_{k+1}, \dots, x_n satisfy ϕ (i.e., $\phi \in \text{EMAJ3SAT}$). so $\sum_{z \in B} T(s, a, z) > 1/2$.

$$\left| \sum_{s \in A} T(s, a, t) - \sum_{r \in B} T(r, a, t) \right| > 1/2.$$

Thus, if $\phi \in \text{EMAJ3SAT}$, then A is not $1/2$ -stable with respect to B . If $\phi \notin \text{EMAJ3SAT}$, then A is $1/2$ -stable with respect to both B and C , and thus the partition into A , B and C is ε -homogeneous. \square

Remark: Notice that the constructed MDP M is a plausible factored MDP. There are two classes of reward states, each with a succinct description in terms of the fluents. The 2TBNs that describe the transitions are also nicely factored. In the 2TBN each node has at most three parents (three synchronous parents for the c_j s, three asynchronous parents for the x_i s, those being x_i , v_0 , and v_1 , and two asynchronous parents for v_0 and v_1). Also no node in the 2TBN has depth greater than 2.

Corollary 2 *The problem of finding an ε -stable partition of a factored MDP with a minimum number of aggregate states is coNP^{PF} -hard.*

Proof. This is shown by the construction for Theorem 1. The constructed MDP M in that proof has a $1/2$ -stable partition with three blocks if and only if the formula $\phi \notin \text{EMAJ3SAT}$. If $\phi \in \text{EMAJ3SAT}$, then the coarsest $1/2$ -stable partition must have more than three blocks. \square

Note that we continue to assume that the partitions can be specified using “small” descriptions. If we want to show that the MINIMUM PARTITION PROBLEM is in a particular class, then the input to the problem must specify not only the MDP and ε , but also bounds on the size of the specification of the partition.

Let us formally define the problem k -PARTITION EXISTENCE PROBLEM as follows:

Instance: A factored MDP and a rational number $0 \leq \varepsilon \leq 1$, an integer k , and a string 1^c .

Question: Is there an ε -homogeneous partition with $\leq k$ blocks that can be specified in $\leq c$ bits?

Proposition 3 *The k -Partition Existence Problem for factored MDPs is in $\exists\text{VCP}$.*

Note that $\text{VCP} = \text{coNP}^{\text{C}} = \text{P}$, and that $\exists\text{VCP}$ is in the Counting Hierarchy. Because of the probabilistic quantifier, it is at least as hard as any of the classes in the Polynomial Hierarchy (i.e., $\text{PH} \subseteq \text{VCP}$), but is still contained in PSPACE.

An Easy Case to Test Stability

There are probably a variety of conditions one can put on the transition functions of the MDP and the type of partitions allowed that would restrict the complexity of the stability-testing problem. Here we outline one such set of restrictions. We do not necessarily consider it a reasonable set, but simply one that brings the complexity of the problem down to a tractable level.

Theorem 4 *The FACTORED MDP ε -HOMOGENEITY problem for a factored MDP is in P if there are constants c_p and c_T such that the following holds:*

1. Each block is always described by a formula on the fluents that mentions only c_p fluents.
2. The probability transition function is represented by a 2TBN containing no synchronous arcs, and at most c_T asynchronous arcs coming into any time $t + 1$ fluent from time t .

Note that Littman showed that a 2TBN with synchronous arcs can be modeled by a similar-sized 2TBN with no synchronous arcs. However, the transformation given in (Littman 1997) increases the in-degrees of those fluents that depended on synchronous fluents, thus affecting the last condition of the Theorem.

Proof. It is sufficient to give a polynomial-time algorithm to determine whether partition block A is ε -stable with respect to partition block B .

Let c_p be the constant bounding the number of fluents in the formula describing a partition block, and let c_T be the constant bounding the number of predecessors of a fluent in a 2TBN. For every action a , we potentially need to calculate

$$\sum_{t \in B} T(s, a, t) \tag{2}$$

for every $s \in A$. For simplicity of notation, assume that the formula ϕ_B describing block B is on fluents x_1, \dots, x_{c_p} . So, for fixed s , the sum given in (2) is calculated from the probability of transitions from s to states satisfying ϕ_B . That probability can be calculated as follows. Consider each of the 2^{c_p} assignments σ to the relevant c_p fluents. (So σ is an assignment to c_p fluents, and a partial assignment to the entire n fluents.) If σ does not satisfy ϕ_B , then discard it. If σ does satisfy ϕ_B , then we calculate the probability of a transition from s to any state with partial assignment σ . That probability depends on the setting of at most $c_T c_p$ variables in state s . The total contribution to the sum (1) from partial assignment σ is equal to $2^{n - c_p}$ times the probability of a transition from s to an arbitrary state obeying partial assignment σ .

So we have shown how to calculate (2) in $O(c_T c_p 2^{c_p})$ time for a fixed $s \in A$. Since c_T and c_p are constant, that is a constant factor. Our real problem is to determine the value of (2), which requires calculating for all $s \in A$. More precisely, we need to calculate the maximum variation in this sum for any two states in A . As above, we need only do separate calculations for the assignments to the fluents mentioned in ϕ_A and their predecessors, so there are really only $2^{c_p c_T}$ partial assignments that need be considered. We can calculate the value of our sum for each one and multiply by the appropriate power of 2 to compute each sum. Since there are only constantly many sums in question, finding the maximum variance is in P. \square

This theorem does not contradict the hardness proofs above. In particular, in the proof of Theorem 1, the factored MDP that is constructed has the property that

even, first depends on a constant number of previous fluents (two bits of v , plus perhaps the three variables corresponding to that clause). However, some of the dependences are synchronous. Furthermore, the partition is defined by a formula on a large number ($m + 2$) of the fluents. Therefore, it does not fit the hypotheses of Theorem 4.

Note that the restrictions in Theorem 4 severely limit both the number of possible partitions for a given 2TBN, and the type of possible 2TBNs. The first restriction seems unreasonable at first glance, in part because it limits the number of possible refinements for each block in a partition. A natural question is whether there are reasonable characterizations of 2TBNs that actually have small ϵ -homogeneous partitions. Is such a "localness" condition (each fluent depending on a very small number of others) either necessary or sufficient for small partitions?

Hardness of Exact Stability Testing

A complete set for $C=P$ is the set of Boolean formulas where *exactly* half the assignments satisfy the formula. The proof of Theorem 5 follows the proof of the analogous theorem for approximate stability testing.

Theorem 5 *The problem. "Is*

$$\left| \sum_{I \in B} T(s, a, t) - \sum_{I \in B} T(s', a, t) \right| = \epsilon?"$$

is complete for the class $C=P$.

Corollary 6 *The exact stability testing problem ($\epsilon = 0$ in Theorem 5) is complete for the class $C=P$.*

One can easily modify the construction that shows that the ϵ -stability question is coNP^{PP} -hard to show that the stability question for exact partitioning is $\text{coNP}^{\text{C=P}}$ -hard. The question remains, what is the class $\text{coNP}^{\text{C=P}}$.

Torán showed (Torán 1988; Torán 1991) that $\text{NP}^{\text{PP}} = \exists \text{CP} = \exists \text{C=P} = \text{NP}^{\text{C=P}}$. He has also shown directly that $\text{coNP}^{\text{PP}} = \text{coNP}^{\text{C=P}}$ (Torán 1999). This shows that the complexity of exact homogeneity testing is exactly the same as the complexity of ϵ -homogeneity testing.

Structured Value Iteration

In the work on BMDPs, the basic idea is to first aggregate states of a 2TBN, and then apply a suitable policy-finding algorithm. Boutilier, Dearden, *et al.*, also construct aggregate MDPs from 2TBNs, but they do so in the course of constructing a policy (Boutilier, Dearden, & Goldszmidt 1995; Boutilier & Dearden 1996; Boutilier, Dearden, & Goldszmidt 1999). The policies constructed are themselves represented as trees. These trees are described by a set of Boolean formulas that partition the state space, and for each element of the

partition, an action. This means that when the specification of a state satisfies a particular formula, the corresponding action is applied.

Remember that 2TBN states are described by assignments of the state variables. For simplicity, they assume that the variables are Boolean. They further assume that the 2TBNs can be described with no synchronous arcs. We can modify the example from the proof of Theorem 1 to meet these criteria and still show that their underlying algorithm must be able to solve NP^{PP} -hard problems. We can remove the synchronous arcs by adding another fluent that controls "time," breaking the process up into two phases, depending on the state of the new variable: assign the x_i s, then set the c_j s. However, the reward criterion depends on a large number of fluents, namely, all the c_j s, violating the first condition of Theorem 4. This complexity does not come as an enormous surprise, since Boutilier, Dearden, and Goldszmidt give an example where the reward function is based on the parity of the binary string identifying the state, and thus is exponentially larger than the number of fluents. They do not claim that their work is a panacea for the curse of dimensionality.

Their algorithms take, as input, a structured policy, and use a form of policy or value iteration to improve the evaluations of each "region" of the state space specified as a leaf of the policy tree, a partial assignment, or a Boolean formula. Based on these values, some or all of the regions may split, and each new region is assigned an action, creating a new tree. In order to compute the policy, another tree, the value tree, is constructed. At each iteration, the value tree is regressed using some variant of value iteration.

The authors then go on to define value trees over value *intervals*, to "collapse" subtrees with similar values down to leaves whose values are represented as intervals. The stability question for these value trees is as computationally complex as is the stability question for Givan *et al.*'s partitions; the construction given in the proof of Theorem 1 carries over almost directly. In order to remove synchronous arcs, we add a new fluent, w , so that when $w = 0$, the action a acts as before on the x_i s but preserves v and sets the c_j s to 0 and w to 1. When $w = 1$, action a sets the x_i s to 0 but sets the c_j according to the values of the appropriate literals at the previous step, sets v as in the previous construction, and flips w . Thus, the previous transitions are stretched over two steps.

The reward function assigns reward 100 to any state with $v = 2$ and all $c_j = 1$ (and thus $w = 0$). Any state with ($v = 0$ or $v = 1$) and $w = 1$ has reward 50; all other states have reward 0.

The initial value tree branches on v . On the $v = 0$ and $v = 1$ branches, it branches on w ; on the $v = 2$ subtree, it branches on each c_j in turn, either to reward 0 or to the next c_j , and finally, if all the $c_j = 1$, to reward 100.

The question is, whether the $v = 0$ and $w = 0$ branch splits when value iteration is applied. This is equivalent

to asking whether the partition is stable, and thus is $C=P$ -hard in the exact case, and $coNP^P$ -hard in the approximate case.

Other Approaches

One way that researchers have further reduced the size of succinctly represented MDPs is a technique called *feature selection*. In such reductions, a domain expert chooses certain features of the state space that are expected to be the primary contributors to the transition probabilities and simply ignores any additional information. Tsitsiklis and Van Roy (Tsitsiklis & Roy 1996) give a careful analysis of dynamic programming algorithms based on feature selection. They do not give any complexity analysis, but they show that such policy-finding algorithms will converge, and will find as good a policy as is possible, given the choice of features. However, this means that their methods still rely on a human or other expert guessing or recognizing the appropriate features, and in fact on there being a few such dominant features. This brings us no closer to a general solution than the existential quantifier ("there exists a partition...") in Givan *et al.*'s work. However, there are definitely cases when feature selection works extremely well.

What our analysis shows is that *testing* a partition of states for ϵ -homogeneity is hard, no matter how the partition was derived. When there is outside knowledge of the system being partitioned, there may be times when the testing can be avoided altogether, or replaced by some sort of sampling (probabilistic verification).

Another approach is to dynamically construct the partition in the course of the policy-finding algorithm. This is a purely heuristic approach, but can lead to significant speed-ups in certain cases, such as those reported in (Boutilier, Dearden, & Goldszmidt 1999; Hoey *et al.* 1999). Naturally, such a heuristic may lead to an approximation of the optimal solution—depending both on the aggregation scheme employed and the policy-finding algorithm—with no guarantee on the quality of the approximation or the running time of the process.

As we mentioned, approximate or exact partition algorithms would potentially yield faster exact or approximate algorithms—implying the collapse of complexity classes we believe or know to be distinct. Therefore, we cannot expect to find a polynomial-time algorithm that yields significantly smaller representations of succinctly represented MDPs.

Therefore, we cannot expect to find a polynomial-time algorithm that yields significantly smaller representations of succinctly represented MDPs.

Acknowledgments

This work partially supported by NSF grants CCR-9610348 and CCR-9800070. This work was done while the first author was visiting The University of Illinois at Chicago. The authors would like to thank Christopher

Lusena and Fred Green for useful conversations, Martin Mundhenk for first observing that the NP-hardness claim could be strengthened, Jacobo Torán for making explicit a result implicit in his dissertation and JACM paper, and Bob Givan for writing up the NP-hardness proof and explaining (more than once!) what they actually assumed in their implementations of the BMDP constructions.

References

- Bellman, R. 1957. *Dynamic Programming*. Princeton University Press.
- Blythe, J. 1999. Decision-theoretic planning. *AI Magazine* 20(2):37–54.
- Boutilier, C., and Dearden, R. 1996. Approximating value trees in structured dynamic programming. In *Thirteenth International Conference on Machine Learning*, 54–62.
- Boutilier, C.; Dean, T.; and Hanks, S. 1999. Decision-theoretic planning: Structural assumptions and computational leverage. *Journal of AI Research* 11:1–94.
- Boutilier, C.; Dearden, R.; and Goldszmidt, M. 1995. Exploiting structure in policy construction. In *14th International Conference on AI*.
- Boutilier, C.; Dearden, R.; and Goldszmidt, M. 1999. Stochastic dynamic programming with factored representations. <http://www.cs.toronto.edu/~cebly/papers.html>.
- Dean, T., and Givan, R. 1997. Model minimization in Markov decision processes. In *Proceedings of the 14th National Conference on Artificial Intelligence and 9th Innovative Applications of Artificial Intelligence Conference (AAAI-97/IAAI-97)*, 106–111. Menlo Park: AAAI Press.
- Dean, T.; Givan, R.; and Leach, S. 1997. Model reduction techniques for computing approximately optimal solutions for Markov decision processes. In Geiger, D., and Shenoy, P. P., eds., *Proceedings of the 13th Conference on Uncertainty in Artificial Intelligence (UAI-97)*, 124–131. San Francisco: Morgan Kaufmann Publishers.
- Fenner, S.; Homer, S.; Green, F.; and Pruim, R. Quantum NP is hard for PH. In *Sixth Italian Conference on Theoretical Computer Science*. World-Scientific.
- Fenner, S.; Homer, S.; Green, F.; and Pruim, R. 1999. Quantum NP is hard for PII. *Royal Society of London A* 455:3953–3966.
- Givan, R., and Dean, T. 1997. Model minimization, regression and propositional STRIPS planning. In *Proceedings of the 15th International Joint Conference on Artificial Intelligence (IJCAI-97)*, 1163–1168. San Francisco: Morgan Kaufmann Publishers.
- Givan, R. L., and Dean, T. 1999. Model minimization for Markov decision processes. Manuscript.
- Givan, R.; Leach, S.; and Dean, T. 1997. Bounded parameter Markov decision processes. In Steel, S., and

- From: AIPS 2000 Proceedings. Copyright © 2000 AAAI (www.aaai.org). All rights reserved.
- Alami, R., eds., *Proceedings of the 4th European Conference on Planning (ECP-97): Recent Advances in AI Planning*, volume 1348 of *LNAI*, 234–246. Berlin: Springer.
- Givan, R. 1999. Personal communication.
- Goldsmith, J., and Mundhenk, M. 1998. Complexity issues in Markov decision processes. In *Proceedings of IEEE Conference on Computational Complexity*, 272–280.
- Goldsmith, J.; Littman, M.; and Mundhenk, M. 1997. The complexity of plan existence and evaluation in probabilistic domains. In *Proc. 13th Conf. on Uncertainty in AI*. Morgan Kaufmann Publishers.
- Harmanec, D. 1999. A generalization of the concept of Markov decision process to imprecise probabilities. In *Proceedings of the 1st International Symposium on Imprecise Probabilities and Their Applications*, 175–182.
- Hoey, J.; St-Aubin, R.; Hu, A.; and Boutilier, C. 1999. SPUDD: Stochastic planning using decision diagrams. In *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, 279–288.
- Littman, M.; Goldsmith, J.; and Mundhenk, M. 1998. The computational complexity of probabilistic plan existence and evaluation. *The Journal of AI Research* 9:1–36.
- Littman, M. L.; Majercik, S. M.; and Pitassi, T. 2000. Stochastic Boolean satisfiability. *Journal of Automated Reasoning*.
- Littman, M. 1997. Probabilistic propositional planning: Representations and complexity. In *Proc. 14th National Conference on AI*. AAAI Press / MIT Press.
- Littman, M. L. 1999. Initial experiments in stochastic satisfiability. In *Proceedings of Sixteenth National Conference on Artificial Intelligence*.
- Lusena, C.; Goldsmith, J.; and Mundhenk, M. 1999. Nonapproximability results for Markov decision processes. In preparation.
- Majercik, M., and Littman, M. 1998a. Maxplan: A new approach to probabilistic planning. In *Artificial Intelligence and Planning Systems*, 86–93.
- Majercik, S. M., and Littman, M. L. 1998b. Using caching to solve larger probabilistic planning problems. In *Proceedings of Fifteenth National Conference on Artificial Intelligence*, 954–959.
- Mundhenk, M.; Goldsmith, J.; Lusena, C.; and Alender, E. 1999. Complexity results for finite-horizon Markov decision process. *Journal of the ACM*. To Appear.
- Papadimitriou, C. 1994. *Computational Complexity*. Addison-Wesley.
- Puterman, M. L., and Shin, M. 1978. Modified policy iteration algorithms for discounted Markov decision problems. *Management Science* 24:1127–1137.
- Toda, S. 1991. PP is as hard as the polynomial-time hierarchy. *SIAM Journal on Computing* 20:865–877.
- Torán, J. 1988. *Structural Properties of the Counting Hierarchies*. Ph.D. Dissertation, Facultat d'Informàtica de Barcelona.
- Torán, J. 1991. Complexity classes defined by counting quantifiers. *Journal of the ACM* 38(3):753–774.
- Torán, J. 1999. Personal communication.
- Tsitsiklis, J. N., and Roy, B. V. 1996. Feature-based methods for large scale dynamic programming. *Machine Learning* 22:59–94.
- Wagner, K. 1986. The complexity of combinatorial problems with succinct input representation. *Acta Informatica* 23:325–356.