# Dynamic decision making in stochastic partially observable medical domains: Ischemic heart disease example.

Milos Hauskrecht

MIT Laboratory for Computer Science, NE43-421

545 Technology Square

Cambridge, MA 02139

*milos@medg.lcs.mit.edu*

**Abstract**

The focus of this paper is the framework of Partially observable Markov decision processes (POMDP) [Astrom 65], [Lovejoy 91a], [Cassandra 94] [Hauskrecht 96]) and its role in modelling and solving complex decision problems in stochastic and partially observable medical domains. The POMDP framework compared to alternative decision making formalisms stresses the dynamic aspect of a decision process, where the quality of any decision is closely tied to and influenced by subsequent decisions and their outcomes, and where states of the world that are relevant for decisions are only partially observable. Moreover the framework allows uniform representation and handling of both investigative actions that enable observations and control actions that induce state changes.

POMDPs were introduced and have been investigated mostly by researchers in control theory and operations research and only recently have become the focus of attention of research in AI, mostly in connection with problems of robot navigation. In the paper I explore the potential of the framework in modelling and solving complex medical decision problems. The selected domain examined is the management of the patient with chronic ischemic heart disease [Wong et.al. 90] [Hauskrecht 96b] .

## Introduction

Over the history of AI in medicine a large amount of research work was devoted to the development of methods and techniques capable of modelling the decision process of a physician in situations in which the effect of actions is uncertain. The focus of the work in this area has gradually shifted from initial simpler problems that were sufficiently modelled by simple lotteries to more complex problems that emphasize the dynamic aspect of the decision process with more decisions affecting and contributing to the final solution. Such problems emerge commonly in situations in which a physician attempts to manage (treat) the patient with some disorder over time while following various objectives, like e.g. the reduction of the length of the treatment or the minimization of the invasiveness of selected procedures etc.

The dynamic decision problem can be described by means of the model of the patient behavior under different interventions and an objective function that quantifies objectives pursued. The task is then to select an action or a sequence of actions such that the objective function is optimized.

This also corresponds to the more general control problem in which one tries to identify an optimal course of actions with regard to the given model of the system and provided objective function.

A typical and frequently used framework for representing dynamic decision problems in stochastic environments is a Markov decision process (MDP) [Howard 60], [Puterman 94]. The framework allows one to describe the stochastic behavior of the controlled system, as well as payoffs associated with various transitions that are in turn exploited in the computation of objective function. However the MDP is based on the assumption that the state of the system (or patient) is always perfectly observable, i.e. that there is no uncertainty associated with what state the system is in once the action is performed. Unfortunately this does not always reflects the nature of many dynamic decision problems in medicine that are characterized by a partial observability of the the underlying patient state. This is seen in the following fairy general scenario:

*The patient suffers from a problem (disorder) that cannot be completely identified based on initial observations. The physician treating the patient tries to correct the problem over time and achieve a set of global objectives (e.g. relieve the patient's symptoms, avoid death etc.) by intervening with various treatment actions and procedures. This process is often complicated by the following: based on the available information the physician is uncertain about the underlying patient state; the patient's response to various actions is not deterministic but it is governed by chance; and the patient's state is not static but evolves over time and can possibly lead to some highly negative states when not treated properly. The physician can reduce his current uncertainty about the patient state by various observations and tests. However they need not allow him to perfectly narrow the underlying patient state, leaving open more diagnostic possibilities. On the other hand some of the more precise tests are risky and carry an additional cost with regard to the problem objectives (some of the investigative procedures can increase the risk of death). Therefore at any point in time the physician faces the problem of choosing the next action/s to do (e.g. treatments, investigative procedures). In considering various options he needs to evaluate the benefits as well as the cost of possible actions from the perspective of his long and short term objectives. This often requires him to examine and compare multiple possible future scenarios.*

A framework that allows us to model both partial observability, as well as control over observations is the Partially observable Markov decision process ([Astrom 65], [Lovejoy 91a], [Cassandra 94], [Hauskrecht 96]). The framework was introduced and has been investigated mostly by researchers in control theory and operations research and only recently it has become the focus of attention of research in AI. The work in AI community has been directed mostly to various navigation problems [Littman et. al. 95a] [Parr, Russell 95] [Hauskrecht 96]. In this paper I will summarize some of the basic features of the POMDP framework, show how one can go about solving decision problems within it, and then examine the potential of the framework to represent the problem of the management of patients with chronic ischemic heart disease.

## Partially observable Markov decision process

A *partially observable Markov decision process (POMDP)* describes the stochastic control process with partially observable states and formally corresponds to 6-tuple $(S, A, \Theta, T, O, C)$ where $S$ is a set of states, $A$ is a set of actions, $\Theta$ is a set of observations, $T$ is a set of transition probabilities between states that describe the dynamic behavior of the modeled environment, $O$ stands for a set of observation probabilities that describe the relationship among observations, states and actions, and $C$ denotes a cost model that assigns costs to state transitions and models payoffs associated with such transitions (alternative formulations can include rewards).
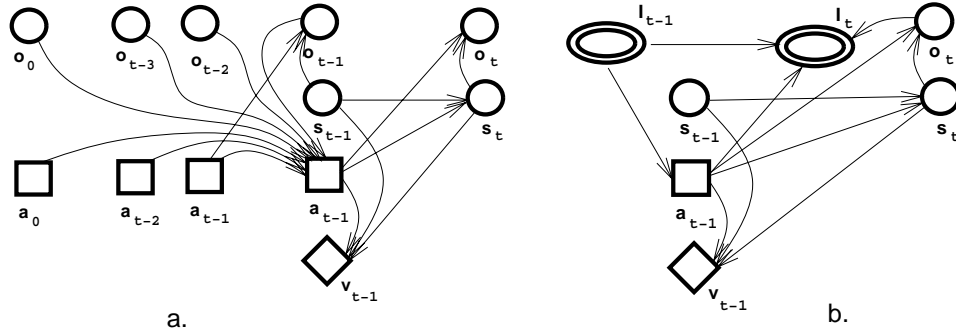
Figure 1: Influence diagrams describing the POMDP.

The *decision (or planning) problem* in the context of POMDP requires one to find an action or a sequence of actions for one or more information states that cause the objective function corresponding to the expected cost incurred over some time horizon to be minimized. An *information state* represents all the information that is available to the agent and that is relevant to the selection of the optimal action. It consists of either a complete history of previous actions and observations or corresponding sufficient statistics. Basic dynamic dependencies between components of the POMDP model can be represented graphically using influence diagrams in figure 1. In the first case the action depends on the complete history while in the second case the information state that stands for either the complete history or its sufficient statistics is used.

A *total expected cost* reflects the quality of the control over time under action outcome uncertainty. It can combine contributions from one step costs of the cost model in different ways. Typical decision criteria for combining one step costs are additive, e.g.:

- finite horizon criterion, that minimizes expected cost for next $n$ steps: $\min E(\sum_{t=0}^{n} c_t)$;

- infinite horizon criteria, that either:

  1. minimizes expected discounted cost $\min E(\sum_{t=0}^{\infty} \gamma^t c_t)$, with $0 \leq \gamma < 1$ being a discount factor, or

  2. minimizes average expected cost per transition $\min \lim_{N \to \infty} \frac{1}{N} \sum_{t=0}^{N} c_t$;

In the following the two of the above criteria: finite horizon and infinite discounted horizon, that are used most often in practice will be assumed.

## Solving the problem

The process of finding the best action (or a sequence of actions) in the POMDP context is closely related to the computation of the objective function that corresponds to the minimum expected cost. For the $n$ step-to-go problem the minimum expected cost can computed using the following

recursive formula that is based on the Bellman's principle of optimality:

$$V_n^*(I_n) = \min_{a \in A} \rho(I_n, a) + \gamma \sum_{o \in \Theta_{next}} P(o|I_n, a) V_{n-1}^*(\tau(I_n, o, a))$$

where $V^*(.)$ is a value function that stands for the minimum expected cost, $I_n$ denotes current information state; $\rho(I_n, a)$ is the expected transition cost from state $I_n$ under action $a$ and can be computed as $\rho(I_n, a) = \sum_{s \in S} \sum_{s' \in S} P(s'|s, a) P(s|I_n) C(s, a, s')$; $\Theta_{next}$ is a set of observations that are available in the next step; $\tau$ is a transition function that maps information state, new action and observation to the next step information state; and $\gamma$ is the possible discount factor. The optimal control for state $I_n$ and the $n$ step-to-go problem then simply corresponds to the action that minimizes the value function, i.e.:

$$\mu_n^*(I_n) = \mathrm{argmin}_{a \in A} \rho(I_n, a) + \gamma \sum_{o \in \Theta_{next}} P(o|I_n, a) V_{n-1}^*(\tau(I_n, o, a))$$

where $\mu(.)$ denotes the optimal control function that maps information states to action space. Identical formulas (less the index denoting the number of steps to go) can be derived for the infinite discounted horizon problem with a stationary control policy.

The problem with finding optimal actions or policies in the POMDP framework is that the task is usually computationally very expensive. This is documented by a fact that even finding the optimal decision for a single initial state and finite horizon is PSPACE-hard [Papadimitriou, Tsitsiklis 87]. This is simply because the number of information states one potentially needs to visit grows exponentially with the number of steps to be explored. A far worse situation emerges when one is required to find the solution for all initial states (so called policy problem). This is also a situation when information state space is properly modeled by a relatively simple belief space for which the value function is known to be finite, piecewise linear and concave [Smallwood, Sondik 73]. The problem here is that the number of linear segments describing the exact value function not only grows exponentially in the size of the action and observation spaces, similar to the problem with a single initial state, but existing algorithms (see e.g. [Smallwood, Sondik 73], [Cheng 88], [Cassandra et.al. 94], [Hauskrecht 96]) for computing linear segments of can use exponential time and space even in the case when the number of useful segments is not exponential. Thus, the price paid for the increased expressivity of the POMDP framework is the computational intractability of the decision/planning problems.

The problem of computational efficiency of exact optimization methods leads naturally to the exploration of various approximation methods and shortcuts that allow one to acquire good solutions with less computation. There are many different approaches one can use for this task and they can be applied to solve both decision and policy problems. Most of these are based on the approximation of value function and use various versions of approximate dynamic programming (for finite horizon case) and approximate value iteration (for infinite discounted horizon). Methods one can employ for this task include (see [Hauskrecht 96]): MDP based approximation, blind policies, point-based approximation methods with either least square error or various interpolation-extrapolation rules, restricted Sondik's method and various methods based on feature extraction mappings. All methods, except feature based, compute solutions by using a finite number of points that sample the information space. On the other hand feature based methods focus primarily on the approximation of the information space.

In the following the attention will be paid to the problem of selecting and improving decision from known approximations of value functions, and various ideas and methods for solving such a problem
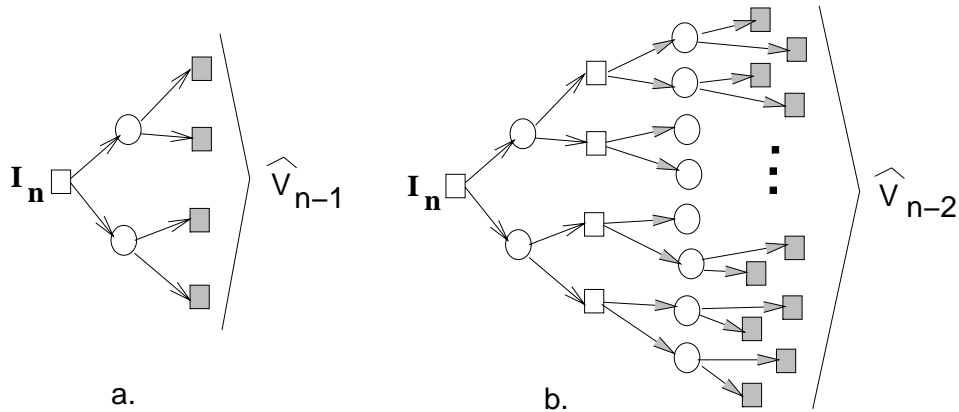
Figure 2: Selecting an action using decision tree expansion.

will be proposed and described. This problem is important also for the purpose of constructing an on-line decision-maker as it is usually easier to store precomputed value function approximation, rather than to store the precomputed policy. The detailed description of various methods that compute value function approximations used in this problem as well as ideas underlying such methods can be found in [Hauskrecht 96].

## Incremental decision algorithms

The easiest way to choose an action is to modify the original optimality formula, by using a value function approximation instead of the optimal value function:

$$\widehat{\mu}_n(I_n) = \operatorname{argmin}_{a \in A} \rho(I_n, a) + \gamma \sum_{o \in \Theta_{next}} P(o|I_n, a)\widehat{V}_{n-1}(\tau(I_n, o, a))$$

This formula and associated computation can be represented using a one step decision tree (see figure 2 a), in which a rectangle corresponds to a decision node that can be associated with an information state and a circle to a chance node that is associated with some information state and a specific action to be performed in that state. In this case the value function approximation is applied to leaves of such a decision tree.

The important thing in this respect is that the same idea of the decision tree expansion can be applied also to the leaves of a one-step decision tree, as it simply corresponds to further unwinding of the recursive formula (see case b in figure 2). By this means one can apply the value function approximation one layer lower, hoping that this would bring the approximation in the first layer closer to the optimal one. Unfortunately this property need not hold in general and is strongly dependent on methods used to compute the value function approximation that is employed. Various methods that satisfy this property were described in [Hauskrecht 96] and they include MDP approximation, blind policy method, restricted Sondik method and some point based approximations. When this property is satisfied for any number of expansions we say that the value function approximation satisfies the *recursive improvement property* [Hauskrecht 96] and the methods above also satisfy this more general property.

The fact that one can improve the value function by expanding the decision tree more can be used in the design of anytime algorithms [Dean 91], that allow to incremental improvement in the value function approximation, i.e. algorithms that can be interrupted anytime and the more time they are given the better is the resulting value function approximation. Such algorithms can be built easily by starting with initial value function approximation and then improving it gradually by a further expansion of the decision tree. If enough time is offered one can even reach optimal solution for some finite horizon problem by simply using the complete decision tree.

The basic incremental improvement algorithm that gradually expands the decision tree, layer by layer can be modified in a number of ways. One possible modification is to change the expansion strategy and allow branches of the decision tree to grow unevenly. Note that given the value function approximation satisfying the recursive improvement property the new approximation with unevenly grown tree branches also guarantees that the property hold. Thus, it is not worse than the previous one. The uneven growth can be exploited whenever it is possible to identify what branches are more likely to improve the approximation more, e.g. by using a suitable heuristic function.

The other important modification of the basic method is based on the capability to prune provably suboptimal branches of the decision tree before they are fully expanded. This can be achieved by using approximations corresponding to value function bounds. Then whenever the value function approximation is a bound of the optimal value function and also satisfies the recursive improvement property then expansion of the tree corresponds to its incremental improvement. Therefore, when value function bounds with such properties are used instead of a single value function approximation, it is possible to incrementally narrow the bound span in which the optimal value function must lay. The capability to compute and narrow value function bounds can be used in a straightforward way to prune some of the suboptimal branches and this can be done often long before they are fully expanded. Note that using incremental improvement with pruning repeatedly drives one to the optimal or near optimal decision (with a guaranteed distance from optimal) and this also for the infinite discounted case. Bounds or bound spans can also be exploited in building heuristical rules to guide the uneven growth of the decision tree [Hauskrecht 96].

Value function approximations corresponding to bounds can be computed using various methods. Relatively simple bounds can be acquired using MDP approximation and the blind policy method (see [Hauskrecht 96]). These two methods are fast and both "ignore" partial observability in different ways. The MDP approximation method uses solutions acquired for a perfectly observable case, i.e.:

$$V_{L,n}(I_n) = \sum_{s \in S} P(s|I_n) V_{MDP,n}^*(s)$$

and computes the lower bound of the optimal value function for the partially observable case. On the other hand the blind policy method uses solutions computed for fixed policies that ignore any information available (thus blind policies) and allows computation of an upper bound. The basic formula is:

$$V_{U,n}(I_n) = \min_{\phi \in \Phi} \sum_{s \in S} P(s|I_n) V_{\phi,n}(s)$$

where $\Phi$ is a set of fixed blind policies and $V_{\phi,n}$ corresponds to the expected cost for a blind policy $\phi$. There are other methods one can use to obtain better bounds for some problems and these can be found in [Lovejoy 91b] [White, Scherer 94], [Hauskrecht 96]. In some cases one can even construct incremental improvement versions of these methods and combine them with incremental methods that expand the decision tree [Hauskrecht 96].

| State variables | Observation variables | Actions |
|---|---|---|
| • **status:**<br>  dead<br>  alive<br>• **coronary artery disease**<br>  normal<br>  mild (no significant stenosis)<br>  moderate (1 or 2 vessels stenosis >70%<br>    no left main coronary artery (LMCA))<br>  severe (3 vessel stenosis, no LMCA)<br>  LMCA stenosis (> 50%)<br>• **ischemia level**<br>  normal<br>  mild<br>  severe<br>  acute MI<br>• **history of MI**<br>  True/False<br>• **history of PTCA**<br>  True/False<br>• **history of CABG**<br>  True/False | • **death**<br>  True/False<br>• **angiogram result**<br>  Positive/Negative<br>• **stress test result**<br>  Positive/Negative<br>• **resting EKG**<br>  Positive/Negative<br>• **chest pain**<br>  no chest pain<br>  typical pain<br>  atypical pain<br>• **acute MI symptoms**<br>  Positive/Negative | no action<br>angiogram investigation<br>stress test<br>medication treatment<br>angioplasty (PTCA)<br>coronary bypass<br>  surgery (CABG) |

Figure 3: Ischemic heart disease - basic model components

There are numerous other methods that can be used to compute a decision. These may include algorithms similar to the presented ones that use value function approximations without the recursive improvement property, that are e.g. computed via various stochastic simulation methods. Although these can occasionaly lead to better actions then those selected by the above methods, the value of a guaranteed improvement or a guaranteed precision of the solution in the case of incremental bound methods is often higher and more important.

# Management of the ischemic heart disease

The expressiveness and capabilities of the POMDP framework can be exploited in representing various problems related to the management or treatment of the patient with diseases that can progress over time. Here the decision about various treatment or investigative options need to be evaluated with regard to the globally pursued goals, including e.g. the patient's well being now and in the future.A typical medical problem that falls into this category is a problem of the management of the chronic ischemic heart disease (IHD) [Wong et.al. 90] [Leong 94] [Hauskrecht 96b]. The objective in this problem is to determine the optimal plan for managing the patient disease with regard to different cost criteria, resulting from the death of the patient or suffering of the myocardial infarction (MI).

The components of the model include state variables, actions and observation variables describing the state of the patient, actions the agent can choose from and observations available to the agent. These are listed together with the corresponding values in the table in figure 3.

The dependencies (both dynamic and causal) between the components of the POMDP model for the IHD decision problem are captured in figure 4. In this figure state variables are represented by
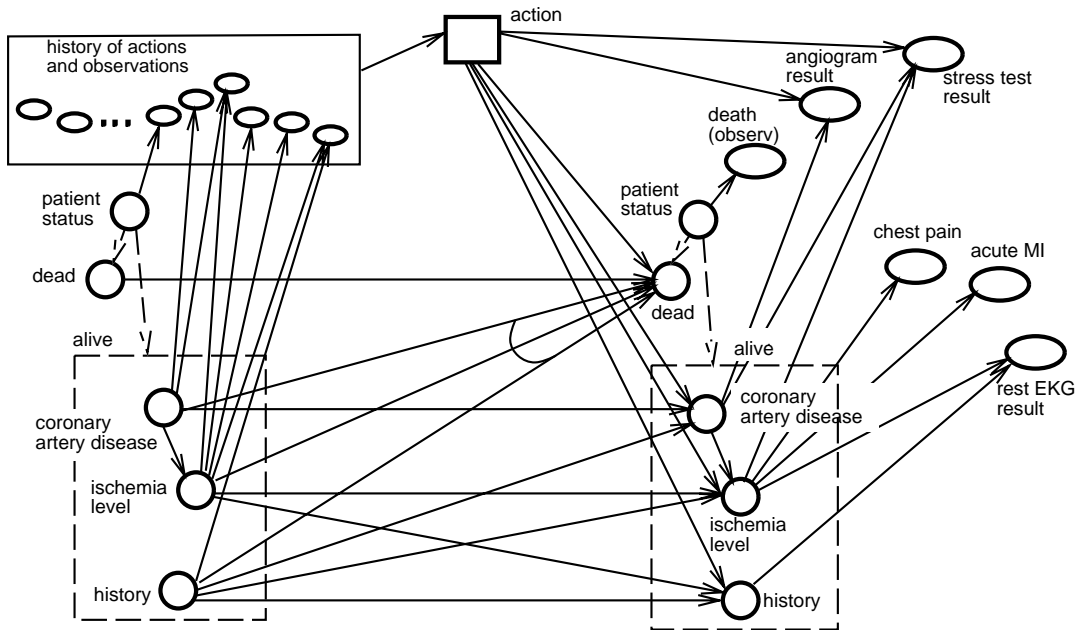
Figure 4: Ischemic heart disease - model of the dynamic behavior

circles, observations as ovals and actions as rectangles. The model is build with the intention to evaluate and reason about the consequences of the long term treatment of the patient with ischemic heart disease. Such a model omits many short term decisions, e.g. one related to handling of the patient with the acute chest pain. This is also the reason why it assumes that the presence of an acute MI can be observed directly.

The *state* of the patient (internal state) in the example is modeled using state variables (represented by circles). A special state variable is associated with global patient status, that differentiate between the two cases: patient being either alive or dead. All other state variables elaborate the case when patient is alive and represent either relevant history information (past MI, or a history of corrective procedures), the current status of coronary arteries, especially from the point of narrowing, and the current level of ischemia.

The set of *actions* in the POMDP model can have exploratory, transitional or cost effects. The exploratory effect of actions is based on their ability to induce observations that in turn can be suggestive of some internal states. An example is an angiogram investigation or stress test. The transition effect of the action is represented by its capability to change the internal state of the patient: e.g. PTCA can lead to the reopening of the blood supply in the main vessels. Note also that actions with intended exploratory effects can lead to changes in the patient state: e.g. increased incidence of MI due to angiogram investigation. The third effect of actions is their cost: the cost measured by the patient's suffering, discomfort and/or economic cost associated with the specific action. Actions that have only an exploratory effect and are neither associated with a cost nor affect the state transition are not explicitly represented in the action set.

*Observations* in the example are modeled through a set of observation variables. These are either triggered through actions, e.g., angiogram result, or corresponds to unconditional observations. Unconditional observations are assumed to be "costless" and available at any time. In our example rest EKG results or chest pain history are assumed to be unconditional due to their relatively low cost. This is unlike the angiogram result that is obtained through risky and costly investigation. In order to focus on the IHD aspect of the disease, we also assume that the acute MI is observed through the auxiliary variable MI-detected.

The *transition* and *observation models* are defined by conditional probability distributions and represent the stochastic nature of the patient's state changes on one side and uncertainty about the actual state on the other. For example the patient with coronary disease can either die, suffer an MI, or get the coronary artery repair as a result of PTCA or CABG, with different probabilities associated with every outcome. Similarly, severe ischemia can, to various degrees, cause typical, atypical or even no chest pain.

The *cost model* describes payoffs associated with possible transitions: for example maximum cost is associated with the transition to the dead state, smaller but still substantial cost is associated with the severe ischemia or occurrence of MI. The decision criteria that try to reduce the expected cost then in fact try to avoid these highly negative states.

The POMDP, like many other frameworks, models continuous time through discretization. In the IHD example it is assumed that every action is associated with a fixed duration and that any change in state occurs between the discretized time points. The way duration of the transitions is set up in the model may in many cases influence the transition probabilities. In the IHD case it is assumed that transitions associated with invasive actions occur within minutes or hours, and transitions associated with non-invasive actions within months (2-3 months). One serious problem with regard to time is the dependence of parameters of the above dynamic model on the age of the patient. That is, there is usually an increased incidence of death for older patients undergoing bypass surgery or PTCA. However the problem with this is that POMDPs and associated solution methods are not good in counting and would then need to consider all possible ages in their model. This problem whenever the problem of long term management is considered can be solved to some extent by including a new state variable named e.g. surgical risk with values denoting various risk groups and probabilistic transitions that tend move patients over time (with increased age) to higher risk groups.

The objective in the context of POMDP is to find an action or a sequence of actions that minimizes the expected cost with regard to the chosen decision criteria. The important thing from the point of action selection is that all known information available to the decision-maker at any point in time can be sufficiently modeled by a belief state that assigns a probability to every possible state. The importance of this stems from the fact that the solution in this case is known to satisfy some nice properties (value function is piecewise linear and concave) that allow one to use better exact or approximation methods. Decision criteria one can use in the IHD case include both finite horizon criterion in which one tries to optimize the treatment with regard to the next $n$ time steps or infinite dicounted horizon criterion that expects one to consider an infinite number of next time steps, with the discounting of more distant future.

## Other possible medical applications of POMDPs

The POMDP model and associated decision criteria provide a useful abstraction that allows a

representation of a fair portion of the mechanism underlying many decision problems conducted in stochastic and partially observable medical domains. Medical problems that fit the above description can also include various emergency room (ER) or intensive care unit (ICU) patient management problems. The typical property of these problems is that the patient state can change dramatically and thus the relevant time discretization should be far shorter than e.g. in the management of the chronic disease. A time criticality often introduces a new dimension of complexity - the need to model the delay with which the observations are made available. This requires to use more complex and computationaly more expensive POMDP models with delayed observations that do not only tradeoff costs and benefits of various observations but also account for their availability in time (see [Hauskrecht 96]. The examples of the time critical patient management problems may include: trauma patient management (see e.g. [Webber et.al. 92] on trauma care), or management of patients with acute chest pain.

## Conclusion

Altough the basic technology for modelling dynamic decision processes via POMDPs has been available for some time it has not been used and exploited in AI in medicine. One reason for this is that the problem solving routines are computationally expensive in the worst case as the decision problem was shown to belong to PSPACE [Papadimitriou, Tsitsiklis 87]. However some approximations that tradeoff the precision for speed can perform well and allow to acquire good solutions efficiently, thus being good alternatives for the optimal methods.

In my research work [Hauskrecht 96] I have described and summarized a number of existing and new problem-solving methods for the POMDP framework that allow one to acquire an optimal or approximate solution faster. In the near future I plan to apply them to the problem of the management of the ischemic heart disease that was described and analysed above. Currently I am in the process of building the underlying POMDP model. I expect that some results will be available in the near future and likely before the submission deadline for the final version of the paper.

## Acknowledgements

## References

[Astrom 65] K.J. Astrom. Optimal control of Markov decision processes with incomplete state estimation. Journal of Mathematical Analysis and Applications, 10,, pp. 174-205, 1965

[Howard 60] R.A. Howard. Dynamic Programming and Markov Processes. MIT press, Cambridge, 1960

[Cassandra et.al. 94] A.R. Cassandra, L.P. Kaebling, M.L. Littman. Acting optimally in partially observable stochastic domains. AAAI-94, pp. 1023-1028, 1994.

[Cassandra 94] A.R. Cassandra. Optimal policies for partially observable Markov decision processes. Brown University, Technical report CS-94-14, 1994.

[Cheng 88] H.-T. Cheng. Algorithms for partially observable Markov decision processes. PhD thesis, University of British Columbia, 1988.

[Dean 91] T. Dean. Decision-theoretic control of inference for time-critical applications. International journal of Intelligent Systems, vol. 6., 1991, pp. 417-441.

[Hauskrecht 96] M. Hauskrecht. Planning and control in stochastic domains with imperfect information. MIT EECS PhD thesis proposal, August 1996, 130 pages.

[Hauskrecht 96b] M. Hauskrecht. Dynamic decision making in stochastic partially observable medical domains. AAAI Spring symposium, pp. 69- 72, 1996.

[Leong 94] T.-Y. Leong. An integrated approach to dynamic decision making under uncertainty. MIT/LCS/TR-631, 1994.

[Lovejoy 91a] W.S. Lovejoy. A survey of algorithmic methods for partially observed Markov decision processes. Annals of Operations Research, 28, pp. 47-66, 1991.

[Lovejoy 91b] W.S. Lovejoy. Computationally feasible bounds for partially observed Markov decision processes. Operations Research, 39:1, pp. 192-175, 1991.

[Littman et. al. 95a] M.L. Littman, A.R. Cassandra, L.P. Kaelbling. Learning policies for partially observable environmnets: scaling up. In Proceedings of the 12-th international conference on Machine Learning, 1995

[Papadimitriou, Tsitsiklis 87] C.H. Papadimitriou, J.N. Tsitsiklis. The complexity of Markov decision processes. Mathematics of Operations Research, 12:3, pp. 441-450, 1987.

[Parr, Russell 95] R. Parr, S. Russell. Approximating optimal policies for partially observable Stochastic domains. IJCAI-95, 1995.

[Puterman 94] M.L. Puterman. Markov Decision Processes: discrete stochastic dynamic programming. John Wiley and Sons, 1994.

[Smallwood, Sondik 73] R.D. Smallwood, E.J. Sondik. The optimal control of Partially observable processes over a finite horizon. Operations Research, 21, pp. 1071-1088.

[Webber et.al. 92] B.L. Webber, R. Rymon, J.R. Clarke. Flexible support for Trauma management through goal directed reasoning and planning. Artificial Intelligence in Medicine, 4:2, 1992.

[White, Scherer 94] C.C. White III, W.T. Scherer. Finite memory suboptimal design for partially observed Markov decision processes. Operations Research, 42:3, pp. 439-455, 1994.

[Wong et.al. 90] J.B. Wong et.al. Myocardial revascularization for chronic stable angina. Annals of Internal Medicine, 113 (1), pp. 852-871, 1990.