

Document downloaded from:

<http://hdl.handle.net/10251/34793>

This paper must be cited as:

García Hernández, MDG.; Ruiz Pinales, J.; Onaindia De La Rivaherrera, E.; Aviña Cervantes, JG.; Ledesma Orozco, S.; Alvarado Mendez, E.; Reyes Ballesteros, A. (2012). New prioritized value iteration for Markov decision processes. *Artificial Intelligence Review*. 37(2):157-167. doi:10.1007/s10462-011-9224-z.



The final publication is available at

<http://link.springer.com/article/10.1007%2Fs10462-011-9224-z>

Copyright Springer Verlag

New Prioritized Value Iteration for Markov Decision Processes

Ma.de Guadalupe Garcia-Hernandez¹, Jose Ruiz-Pinales¹, Eva Onaindia², Alberto Reyes-Ballesteros³, J. Gabriel Aviña-Cervantes¹, Sergio Ledesma-Orozco¹, Edgar Alvarado-Mendez¹

¹ *University of Guanajuato, Comunidad de Palo Blanco s/n, Salamanca, Guanajuato, Mexico, tel. 52 464 6479940
{garciag,pinales,avina,selo,ealvarad}@salamanca.ugto.mx*

² *Universitat Politècnica de València, DSIC, Camino de Vera s/n, 46022, Valencia, España, tel. 34 963877000, onaindia@dsic.upv.es*

³ *Electrical Research Institute, Reforma 113, 62490, Temixco, Morelos, Mexico, tel. 52 7773623811, areyes@iee.org.mx*

Abstract: The problem of solving large Markov decision processes accurately and quickly is challenging. Since the computational effort incurred is considerable, current research focuses on finding superior acceleration techniques. For instance, the convergence properties of current solution methods depend, to a great extent, on the order of backup operations. On one hand, algorithms such as topological sorting are able to find good orderings but their overhead is usually high. On the other hand, shortest path methods, such as Dijkstra's algorithm which is based on priority queues, have been applied successfully to the solution of deterministic shortest-path Markov decision processes. Here, we propose an improved value iteration algorithm based on Dijkstra's algorithm for solving shortest path Markov decision processes. The experimental results on a stochastic shortest-path problem show the feasibility of our approach.

Keywords: *Markov decision processes, priority queues, Dijkstra's algorithm.*

1. Introduction

In planning under uncertainty, the planner's objective is to find a policy that optimizes some expected utility. Most approaches for finding such policies are based on decision-theoretic planning (Boutilier 1999) (Bellman 1954) (Puterman 1994). Among these, Markov decision processes (MDPs) constitute a mathematical framework for modeling and deriving optimal policies. Value iteration is a dynamic programming algorithm (Bellman 1957) for solving MDPs, but it is usually not considered because of its slow convergence (Littman 1995). This is because its speed of convergence depends strongly on the order of the computations (or backups).

The slow convergence of value iteration for solving large MDPs is usually tackled up by using one of two approaches (Dai 2007a): heuristic search (Hansen 2001) (Bhuma 2003) (Bonet 2003a,b, 2006), or prioritization (Moore 1993) (Ferguson 2004) (Dai 2007b) (Wingate 2005). In the first case, heuristic search (combined with dynamic programming) is used to reduce the number of relevant states as well as the number of search expansions. Hansen *et al.* (Hansen 2001) considered only part of the state space by constructing a partial solution graph, searching implicitly from the initial state towards the goal state, and expanding the most promising branch of an MDP according to a heuristic function. Bhuma *et*

al. (Bhuma 2003) extended this approach by using a bidirectional heuristic search algorithm. Bonet *et al.* (Bonet 2003a,b) proposed two other heuristic algorithms that use a clever labeling technique to mark irrelevant states. Later on, they explored depth-first search for the solution of MDPs (Bonet 2006). In the second case, prioritization methods are based on the observation that, in each iteration, the value function usually changes only for a reduced set of states. So, they prioritize each backup in order to reduce the number of evaluations (Moore 1993) (Dai 2007b). Ferguson *et al.* (Ferguson 2004) proposed another prioritization method called focused dynamic programming, where priorities are calculated in a different way than in prioritized sweeping. Dai *et al.* (Dai 2007a) extended Bhuma *et al.*'s idea (Bhuma 2003) by using concurrently different starting points. In addition, they also proposed (Dai 2007b) a topological value iteration algorithm, which groups states that are mutually and causally related together in a meta-state for the case of strongly connected states (or MDPs with cyclic graphs). Likewise, other approaches such as topological sorting (Wingate 2005) and shortest path methods (McMahan 2005a,b) have been proposed. On the first hand, topological sorting algorithms can be used to find good backup orderings but their computational cost is usually high (Wingate 2005). On the other hand, shortest path methods have been applied to the solution of MDPs with some success (McMahan 2005a,b).

We consider the problem of finding an optimal policy in a class of positive MDPs with absorbing terminal states, which are equivalent to stochastic-shortest-path problems (Bertsekas 1995). McMahan *et al.* (McMahan 2005a) proposed a method called improved prioritized sweeping (IPS) for solving single goal stochastic-shortest-path problems based on the Dijkstra's algorithm. The advantages of this method are the reduction to Dijkstra's algorithm for the case of Markov chains (or acyclic deterministic MDPs), and the improvement in speed when compared with other methods such as prioritized sweeping (Moore 1993) and focused dynamic programming (Dai 2007b). Unfortunately, IPS has no guaranteed convergence to the optimal policy for the case of stochastic shortest path MDPs (Dai 2007a,b) (Li 2009). Thus, in this work we propose a new prioritized value iteration algorithm based on Dijkstra's algorithm which has guaranteed convergence for the case of stochastic-shortest-path problems in addition that it can deal with multiple goal and start states.

This paper is organized as follows: first we present a brief introduction to MDPs as well as solution methods, then we describe the prioritized sweeping approaches, after that, we describe our algorithm and present experimental results. Finally, we present the conclusions.

2. Markov Decision Processes

Markov decision processes (MDPs) provide a mathematical framework for modeling sequential decision problems in uncertain dynamic environments (Bellman 1957) (Puterman 2005).

Formally, a MDP is a four-tuple (\mathcal{S}, A, P, R) , where \mathcal{S} is a finite set of states $\{s_1, \dots, s_n\}$, A is a finite set of actions $\{a_1, \dots, a_n\}$, $P : \mathcal{S} \times A \times \mathcal{S} \rightarrow [0, 1]$ is the transition probability function, which associates a set of probable next states to a

given action in the current state. $P(a, s, s')$ denotes the transition-probability to reach state s' , if one applies action a in state s . $R(s, a)$ denotes the reward obtained if one applies action a in state s . $\rho(s)$ denotes a policy (or strategy), it yields an action for each state, it is a rule that specifies which action should be taken in each state. The *Markovian property* guarantees that s' only depends on the pair (s, a) . The core problem of MDPs is to find the optimal policy to maximize the expected total reward (Puterman 2005). The value function, which is the expected reward (or utility) when starting at state s and following policy ρ , is given by:

$$V^\rho(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \rho(s_t)) \mid s_0 = s \right] \quad (1)$$

where $\gamma \in [0, 1]$ is a discount factor, which may be used for decreasing exponentially future rewards. For the case of discounted MDPs ($0 < \gamma < 1$), the utility of an infinite state sequence is always finite. So, the discount factor expresses that future rewards have less value than current rewards (Russell 2004). For the case of additive MDPs ($\gamma = 1$) and infinite horizon, the expected total reward may be infinite and the agent must be guaranteed to end up in a terminal state.

Let $V^*(s)$ be the optimal value function given by:

$$V^*(s) = \max_{\rho} V^\rho(s). \quad (2)$$

The optimal value function satisfies the Bellman equation (Bellman 1954) (Puterman 2005) that is given by:

$$V_t^*(s) = \max_a \left[R(s, a) + \gamma \sum_{s'} P(a, s, s') V_{t+1}^*(s') \right]. \quad (3)$$

Value iteration, policy iteration and linear programming are three of the most well known techniques for finding the optimal value function $V^*(s)$ and the optimal policy ρ^* for infinite horizon problems (Chang 2007). However, policy iteration and linear programming are computationally expensive techniques when dealing with problems with large state spaces because they both require solving several linear systems (of equations) of the same size as the state space. In contrast, value iteration avoids this problem by using a recursive approach typically used in dynamic programming (Chang 2007).

Starting from an initial value function, value iteration applies successive updates to the value function for each $s \in S$ by using:

$$\hat{V}(s) = \max_a \left[R(s, a) + \gamma \sum_{s'} P(a, s, s') \hat{V}(s') \right]. \quad (4)$$

Let $\{V_n \mid n = 0, 1, K\}$ be the sequence of value functions obtained by value iteration. Then, it can be shown that every value function obtained by value iteration satisfies $|V_n - V^*| \leq \gamma^n |V_0 - V^*|$. Thus, from the Banach fixed point

theorem, it can be inferred that value iteration converges to the optimal value function $V^*(s)$. One advantage of value iteration comes from the fact that the value functions obtained can be used as bounds for the optimal value function (Tijms 2003).

The computational complexity of one update of value iteration is $O(|S|^2|A|)$. However, the number of required iterations can be very large. Fortunately, it has been shown in (Littman 1995) that an upper bound for the number of iterations required by value iteration to reach an ϵ -optimal solution is given by:

$$n_{it} \leq \frac{b + \log(\frac{1}{\epsilon}) + \log(\frac{1}{1-g}) + 1}{1-g} \quad (5)$$

where $0 < g < 1$, b is the number of bits used to encode rewards and state transition probabilities, and ϵ is the threshold of the *Bellman error* (Puterman 2005) given by:

$$B^t(s) = \max_{a \in A} \left[R(s,a) + g \sum_{s' \in S} P(a,s,s') V^t(s') - V^t(s) \right] \quad (6)$$

The convergence of value iteration may be quite slow for g close to one. For this reason, several improvements to value iteration have been proposed (Puterman 2005). For instance, common techniques may improve convergence rate, reduce the time taken per iteration and/or use better stopping criteria.

One of the easiest ways to improve convergence rate is to update the value functions as soon as they become available (also known as asynchronous updates). For instance, Gauss-Seidel value iteration uses the following update equation (Puterman 2005):

$$V^t(s) = \max_a \left[R(s,a) + g \sum_{s' \in S} P(a,s,s') V^t(s') + g \sum_{s' \in S} P(a,s,s') V^{t-1}(s') \right] \quad (7)$$

It is well known that policy iteration converges in less number of iterations than value iteration does, but it is more expensive per iteration because it requires solving a system of linear equations at each one of the iterations. In contrast, value iteration does not require the solution of any linear system of equations. A combined approach (modified policy iteration) can exploit the advantages of both. Thus, modified policy iteration uses a partial policy evaluation step based on value iteration (Puterman 2005).

Other way of improving the convergence rate as well as the iteration time is using prioritization and partitioning (Wingate 2005). Generally, prioritization methods are based on the observation that, at each iteration, the value function usually changes only for a reduced set of states. Thus, by restricting the computation to only those states, a reduction of the iteration time is expected. It has been outlined that for acyclic problems the ordering of the states, where the

transition matrix becomes triangular may result in a significant reduction in time (Wingate 2005).

Another method to reduce the iteration time is to identify and eliminate suboptimal actions (Puterman 2005). For instance, bounds of the optimal value function can be used to eliminate suboptimal actions. The advantage of this approach is that the action set is progressively reduced with the consequent reduction in time.

On the other hand, the number of iterations can be slightly reduced by using improved stopping criteria based on tighter bounds of the Bellman error (see Eq. (6)) (Puterman 2005). For instance, a stopping criterion would be to stop value iteration when the span of the Bellman error falls below a certain threshold (Puterman 2005).

Last, for the case of large MDPs with sparse transition matrices, memory savings can be obtained by using a *sparse* representation (Agrawal 2002) where only non-zero transition probabilities are stored. In this way, it is possible to handle larger problems than the ones that can be solved otherwise (mainly in highly sparse MDPs). For instance, an adjacency list containing all the state transitions with non-zero probability can be built.

3. Priority-based methods for solving MDPs

Although value iteration is a powerful algorithm for solving MDPs, it has some potential problems. First, some backups are useless because not all states change in a given iteration (Dai 2007b). Second, backups are not performed in an optimal order. Priority-based methods such as prioritized sweeping (PS) (Moore 1993) avoid these problems by ordering and performing backups so as to perform the least number of backups (Dai 2007b). To be more precise, PS maintains a priority queue for ordering backups intelligently. This priority queue is updated as the algorithm sweeps through the state space. PS can begin by inserting the goal state in the priority queue when it is used in an offline dynamic programming algorithm, such as value iteration. At each step, PS pops a state s from the queue with the highest priority and performs a Bellman backup of that state. If the Bellman residual of state s is greater than some threshold value ε or if s is the goal state, then PS inserts its predecessors into the queue according to their priority (Dai 2007b). Unfortunately, the use of a priority queue for all the states of the model may result in an excessive overhead for real-world problems (Wingate 2005), especially for cyclic MDPs.

Focused dynamic programming (Ferguson 2004) is another variant of prioritized sweeping that exploits the knowledge of the start state to focus its computation on states that are reachable from that state. To do this, focused dynamic programming uses a priority metric that is defined using two heuristic functions: an admissible estimate of the expected cost for reaching the current state from the start state and an estimate of the expected cost for reaching the goal state from the current state. In contrast to other forms of prioritized sweeping, this approach removes the state with the lowest priority value from the priority queue,

instead of removing the state with the highest priority value, since it is interested in states through which the shortest path passes.

Dibangoye *et al.* (Dibangoye 2008) proposed an improved topological value iteration algorithm (iTVI) which uses a static backup order. Instead of minimizing the number of backups per iteration or eliminating useless updates, this algorithm attempts to minimize the number of iterations by using a good backup order (a topological order). First, depth-first-search is used to collect all reachable states from the start state. Next, breadth-first-search is used to build a metric $d(s)$, which is defined as the distance from the start state to state s . A static backup order is built from the resulting metric in such a way that states that are closer to the start state be updated first. The algorithm is guaranteed to converge to the optimal value function because it updates all states recursively in the same way as value iteration does.

Meuleau *et al.* (Meuleau 2006) solved stochastic over-subscription planning problems (SOSPs) by means of a two-level hierarchical model. They exploit this hierarchy by solving a number of smaller factored MDPs. Shani *et al.* (Shani 2008) extended the use prioritization to partially observable Markov Decision processes. In this case, backups are prioritized by using the Bellman error as a priority metric and no priority queue is used.

In contrast with the above methods, it is worth to mention a prioritization method that does not require a priority queue (Dai 2007c), instead, it uses a FIFO (first input, first output) queue if the backwards traversal of the policy graph is breadth-first (forwards value iteration), or a LIFO (last input, first output) queue if the backwards traversal is depth-first (backwards value iteration). In both cases, unnecessary backups can be avoided by using a labeling technique (Bonet 2003a,b) and the decomposition of the state space into a number of strongly connected components. Unfortunately, it has been shown that the backup order induced by these algorithms is not optimal (Dai 2007c).

Since the performance of PS depends on the priority metric that it is used to order states in the priority queue, several researchers have investigated alternative priority metrics. For instance, IPS (McMahan 2005a,b) uses a combination of priority metrics (a value change metric, and an upper bound metric). In fact, it has been shown that IPS may outperform other prioritized sweeping algorithms (Dai 2007a,b).

4. Proposed Algorithm

Dijkstra's algorithm is an efficient greedy algorithm for solving the single-source shortest path problem in a weighted acyclic graph. This algorithm is a special case of the A* algorithm but unlike the A* algorithm, Dijkstra's algorithm is not goal oriented. This is because Dijkstra's algorithm computes all shortest paths from a single source node to all nodes and thus solves the one-to-all shortest path problem. In fact it has been shown that Dijkstra's algorithm is a successive approximation method to solve the dynamic programming equation for the shortest path problem (Sniedovich 2006, 2010) and therefore it is based on the Bellman's optimality principle. The main difference between Dijkstra's algorithm

and other dynamic programming methods for the shortest path problem is the particular order in which it processes states; it processes states according to a greedy best first rule. More precisely, Dijkstra's algorithm chooses the next state to be processed as the one having the smallest value of the dynamic programming functional. One implication for the solution of MDPs is that, instead of a topological order, a more suitable update order may be to choose the next state to be updated as the one having the highest value function. For that reason, in our algorithm we use the current value function as a priority metric.

One of the advantages of value iteration and its variants (in particular PS) is that their convergence to the optimal value function is guaranteed for the case of discounted MDPs and for the case of additive MDPs with absorbing states (Bertsekas 1995) (Li 2009). This is because successive applications of the Bellman equation guarantee convergence to the optimal value function. For that reason, in our algorithm (Improved Prioritized Value Iteration, IPVI) we update all predecessors of the best state (having the highest value function) by using the Bellman equation.

Let $V^t(s)$ be the expected cost at time t to reach a goal state starting from a state $s \in S$, $\pi^t(s)$ be the best policy or action at time t and state s , $R(s, a)$ be the reward for action $a \in A$ in state s , $L = \{(s_k, s'_k, a_k, p_k) \mid p_k = P(s'_k \mid s_k, a_k) \neq 0\}$ be the set of all the possible state transitions, G be the set of goal states, γ be the discount factor and ϵ be the maximum error. Basically, our method performs an initialization step followed by successive prioritized removals of each state in the queue with an update of its predecessors by using the Bellman equation until the queue is empty.

As shown in Algorithm 1, for each state $s \in S$, we make $\pi^0(s) = -1$, then if $s \notin G$ then we assign a very large positive constant M to $V^0(s)$ otherwise we set its value to zero. Next, we push each goal state $s \in G$ into the priority queue according to its priority $V^0(s)$. Then, we repeat the following procedure until the priority queue is empty. We pop the state s with the highest priority out from the queue, and then, we update all its predecessors $y \in \text{pred}(s)$. For every update of a predecessor $y \in \text{pred}(s)$ of state s , we compute the Bellman equation

$$V^{t+1}(y) = \max_a \left\{ R(y, a) + \gamma \sum_{\forall (s_k = y, s'_k, a_k = a, p_k) \in L} p_k V^t(s'_k) \right\}, \quad (8)$$

and if $|V^{t+1}(y) - V^t(y)| > \epsilon$ then we push state y into the queue according to its priority $V^{t+1}(y)$, if state y is already in the queue, then its priority is only updated.

[Insert Algorithm 1 about here]

5. Description of Experiments

For the validation of the proposed algorithm, we chose the sailing strategies problem (Vanderbei 1996, 2008), which is a finite state-action-space stochastic-shortest-path problem, in which a sailboat has to find the shortest path between two points of a lake under fluctuating wind conditions.

The details of the problem are as follows: the sailboat's position is represented as a pair of coordinates on a grid of finite size. The sailor has eight actions giving the direction to a neighboring grid position. Each action has a cost (required time) depending on the direction of the boat's heading and the wind. For the action whose direction is just the opposite of the direction of the wind, the respective cost must be high. For example, if the wind is at 45 degrees measured from the boat's heading (*upwind* tack), it requires four seconds to sail from one waypoint to one of the nearest neighbors. But, if the wind is at 90 degrees from the boat's heading (*crosswind* tack), the boat moves faster through the water and can reach the next waypoint in only three seconds. If the wind is a quartering tailwind (*downwind* tack), it requires just two seconds. Finally, if the boat is sailing directly downwind (*away* tack), it requires only one second. Otherwise, the wind can hit the left or right side of the boat (a *port* or a *starboard* tack, respectively). When changing from a port to a starboard tack (or vice versa), it wastes three seconds (*delay*) for every such change of tack. To keep our model simple, we assume that the wind intensity is constant but its direction may change at any time. The wind could come from one of three directions: either from the same direction as the old wind or from 45 degrees to the left or to the right of the old wind. Table 1 shows the probabilities of a change on the wind direction used in all the experiments. When the heading is along one of the diagonal directions, the time is multiplied by $\sqrt{2}$ to account for the somewhat longer distance that must be traveled. Each state \mathcal{S} of the MDP corresponds to a position of the boat (x, y) , a tack $t \in \{0, 1, 2\}$ and a wind direction $w \in \{0, 1, \dots, 7\}$.

[Insert Table 1 about here]

All the experiments were performed on a 2.66 GHz Pentium D computer with 2 GB RAM running Windows XP. All the tested algorithms were implemented using the Java language. The initial and maximum size of the stack of the Java virtual machine was set to 1024 MB and 1536 MB, respectively. For all the experiments, we set $\varepsilon = 10^{-7}$ and $\gamma = 1$. So, we are dealing with an additive MDP, where convergence is not guaranteed by the Banach fixed point theorem (Blackwell 1965). Fortunately, the presence of absorbing states (states with zero reward and 100% probability of staying in the same state) may allow the algorithm to converge (Hinderer 2003). The lake size was varied from 50×50 to 260×260 and the resulting number of states varied from 55296 to 1597536, respectively (without the bounding beaches). We repeated each run 10 times and then we calculated the mean and standard deviation of the solution time.

6. Experimental Results

We tested our approach (IPVI) and several variants of value iteration including different acceleration techniques (Puterman 2005): Gauss-Seidel Value Iteration (GSVI); Gauss-Seidel Value Iteration with updates of only those states (as well as their neighbors) whose value function changed in the previous iteration (GSVI2); and Gauss-Seidel Value Iteration with the same acceleration procedures as GSVI2 plus static reordering of the states in decreasing order of maximum reward (GSVI3), Policy Iteration (PI), and Modified Policy Iteration (MPI). Also other two algorithms were tested: a dynamic programming approach (VDP) (Vanderbei 1996, 2008) and the improved topological value iteration (iTVI) (Dibangoye 2008). We also tested iTVI with different priority metrics but it converged to the optimal policy only for the priority metric suggested by Dijkstra's algorithm.

Figure 1 shows the solution time for all tested algorithms as a function of the number of states. As we can see, IPVI yielded the lowest solution time whereas VDP yielded the highest solution time. For instance, for 525696 states, our algorithm took 29.9 seconds, whereas GSVI3 took 173.6 seconds, GSVI2 took 203.3 seconds, GSVI took 303.6 seconds and VDP took 589.5 seconds. In this case, our algorithm was 5.8 times faster than GSVI3, 6.8 times faster than GSVI2, 10.2 times faster than GSVI, and 19.7 times faster than VDP. For 940896 states, our algorithm took 57.4 seconds, whereas GSVI3 took 412 seconds, GSVI2 took 486.5 seconds, GSVI took 755.1 seconds and VDP took 1376.6 seconds. In this case, our algorithm was 7.2 times faster than GSVI3, 8.5 times faster than GSVI2, 13.2 times faster than GSVI, and 24 times faster than VDP. For 1359456 states, our algorithm took 81.5 seconds. As we can see, iTVI (Dibangoye 2008) was not tested for more than 400000 states because it exhausted the memory resources.

[Insert Figure 1 about here]

Figure 2 shows a closer look of the solution time as a function of the number of states for IPVI, VI with asynchronous updates (GSVI) and iTVI. For instance, for 393216 states, our algorithm took 22.2 seconds, whereas GSVI took 195.4 seconds and iTVI took 347.6 seconds. In this case, our algorithm was 8.8 times faster than GSVI and 15.7 times faster than iTVI.

[Insert Figure 2 about here]

7. Conclusions

In this paper we have proposed and tested a new prioritized value iteration algorithm based on Dijkstra's algorithm for solving stochastic shortest path MDPs. Unlike other prioritized approaches such as IPS, our approach can deal with multiple start and goal states, and since it successively updates each state by using the Bellman equation, it has guaranteed convergence to the optimal solution. In addition, our algorithm uses the current value function as a priority metric since Dijkstra's algorithm suggests that a more suitable update order is given by the value of the dynamic programming functional.

We compared the performance of our method with other state-of-the-art algorithms including different acceleration techniques. At least in the sailing strategies problem, our approach was the fastest algorithm in addition that it has guaranteed convergence to the optimal value function.

References

- Agrawal, S. and Roth, D. (2002). Learning a Sparse Representation for Object Detection. Proceedings of the 7th European Conference on Computer Vision, pp 1-15, Copenhagen, Denmark.
- Bellman, R. E. (1954). The Theory of Dynamic Programming. Bull. Amer. Math. Soc., 60: 503-516.
- Bellman, R. E. (1957). Dynamic Programming. Princeton University Press, N. J., USA.
- Bertsekas, D. P. (1995). Dynamic Programming and Optimal Control. Athena Scientific, Massachusetts, USA.
- Bhuma, K. and Goldsmith, J. (2003). Bidirectional LAO* Algorithm. Proceedings of Indian International Conferences on Artificial Intelligence, pp 980-992.
- Blackwell, D. (1965). Discounted dynamic programming. Annals of Mathematical Statistics, 36: 226-235.
- Bonet, B. and Geffner, H. (2003a). Faster Heuristic Search Algorithms for Planning with uncertainty and full feedback. Proceedings of the 18th International Joint Conference on Artificial Intelligence, Morgan Kaufmann, Acapulco, México, pp 1233-1238.
- Bonet, B. and Geffner, H. (2003b). Labeled RTDP: Improving the Convergence of Real-Time Dynamic Programming. Proceedings of the International Conference on Automated Planning and Scheduling, pp 12-21, Trento, Italy.
- Bonet, B. and Geffner, H. (2006). Learning depth-first search: A unified approach to heuristic search in deterministic and non-deterministic settings and its application to MDP. Proceedings of the 16th International Conference on Automated Planning and Scheduling, Cumbria, UK.
- Boutilier, C., Dean, T. and Hanks, S. (1999). Decision-Theoretic Planning: Structural Assumptions and Computational Leverage. Journal of Artificial Intelligence Research, 11: 1-94.
- Chang, I. and Soo, H. (2007). Simulation-based algorithms for Markov decision processes. Communications and Control Engineering, Springer Verlag London Limited.
- Dai, P. and Goldsmith, J. (2007a). Faster Dynamic Programming for Markov Decision Processes. Technical Report, Doctoral Consortium, Department of Computer Science and Engineering, University of Washington.
- Dai, P. and Goldsmith, J. (2007b). Topological Value Iteration Algorithm for Markov Decision Processes. Proceedings of the 20th International Joint Conference on Artificial Intelligence, pp 1860-1865, Hyderabad, India.
- Dai, P. and Hansen, E. A. (2007c). Prioritizing Bellman Backups without a Priority Queue. Proceedings of the 17th International Conference on Automated Planning and Scheduling, Association for the Advancement of Artificial Intelligence, pp 113-119, Rhode Island, USA.
- Dibangoye, J. S., Chaib-draa, B. and Mouaddib A. (2008). A Novel Prioritization Technique for Solving Markov Decision Processes. Proceedings of the 21st International FLAIRS (The Florida Artificial Intelligence Research Society) Conference, Association for the Advancement of Artificial Intelligence, Florida, USA.

- Ferguson, D. and Stentz, A. (2004). Focused Propagation of MDPs for Path Planning. Proceedings of the 16th IEEE International Conference on Tools with Artificial Intelligence, pp 310-317.
- Hansen, E. A. and Zilberstein, S. (2001). LAO: A Heuristic Search Algorithm that finds solutions with Loops. Artificial Intelligence, 129: 35-62.
- Hinderer, K. and Waldmann, K. H. (2003). The critical discount factor for Finite Markovian Decision Processes with an absorbing set. Mathematical Methods of Operations Research, Springer Verlag, 57: 1-19.
- Li, L. (2009). A Unifying Framework for Computational Reinforcement Learning Theory. PhD Thesis, The State University of New Jersey, New Brunswick, NJ, USA, October.
- Littman, M. L. and Dean, T. L. and Kaelbling, L. P. (1995). On the Complexity of Solving Markov Decision Problems. Proceedings of the 11th International Conference on Uncertainty in Artificial Intelligence, pp 394-402, Montreal, Quebec.
- Meuleau, N., Brafman, R. and Benazera, E. (2006). Stochastic Over-subscription Planning using Hierarchies of MDPs. Proceedings of the 16th International Conference on Automated Planning and Scheduling, pp 121-130, Cumbria, UK.
- McMahan, H. B. and Gordon, G. (2005a). Fast Exact Planning in Markov Decision Processes. Proceedings of the 15th International Conference on Automated Planning and Scheduling, Monterey, CA, USA.
- McMahan, H. B. and Gordon, G. (2005b). Generalizing Dijkstra's Algorithm and Gaussian Elimination for Solving MDPs. Technical Report, Carnegie Mellon University, Pittsburgh, PA, USA.
- Moore, A. and Atkeson, C. (1993). Prioritized Sweeping: Reinforcement Learning with less data and less real time. Machine Learning, 13: 103-130.
- Puterman, M. L. (1994). Markov Decision Processes. Wiley Editors, New York, USA.
- Puterman, M. L. (2005). Markov Decision Processes. Wiley Inter Science Editors, New York, USA.
- Russell, S. (2004). Artificial Intelligence: A Modern Approach. 2nd Edition, Making Complex Decisions (Ch-17), Pearson Prentice Hill Ed., USA.
- Shani, G., Brafman, R. and Shimony, S. (2008). Prioritizing Point-based POMDP Solvers. IEEE Transactions on Systems, Man. and Cybernetics, Vol. 38, No. 6, pp 1592-1605 December.
- Sniedovich, M. (2006). Dijkstra's algorithm revisited: the dynamic programming connexion. Control and Cybernetics, Vol.35, pp 599-620.
- Sniedovich, M. (2010). Dynamic Programming: Foundations and Principles. Second Edition, Pure and Applied Mathematics Series, Taylor and Francis Publishers.
- Tijms, H. C. (2003). A First Course in Stochastic Models. Wiley Ed., Discrete-Time Markov Decision Processes (Ch-6), UK.
- Vanderbei, R. J. (1996). Optimal Sailing Strategies. Statistics and Operations Research Program, University of Princeton, USA (<http://orfe.princeton.edu/~rvdb/sail/sail.html>).
- Vanderbei, R. J. (2008). Linear Programming: Foundations and Extensions. Springer Verlag, 3rd Edition, January.
- Wingate, D. and Seppi, K. D. (2005). Prioritization Methods for Accelerating MDP Solvers. Journal of Machine Learning Research, 6: 851-881.

Algorithm 1 Improved Prioritized Value Iteration (IPVI).

```
IPVI( $R, L, S, G, g, \epsilon$ )
 $(\forall s \in S)V(s) \leftarrow M$ 
 $(\forall s \in G)V(s) \leftarrow 0$ 
 $(\forall s \in G)\text{queue.enqueue}(s, V(s))$ 
while ( $\neg\text{queue.isempty}()$ )
   $s \leftarrow \text{queue.pop}()$ 
  for all  $y \in \text{pred}(s)$ 
     $V'(y) \leftarrow V(y)$ 
    
$$V(y) = \max_a \left\{ R(y, a) + \gamma \sum_{\forall (s_k=y, s'_k=a_k=a, p_k) \in L} p_k V'(s'_k) \right\}$$

    if  $|V(y) - V'(y)| > \epsilon$  then
       $\text{queue.decreasepriority}(y, V(y))$ 
    end
  end
end
return
```

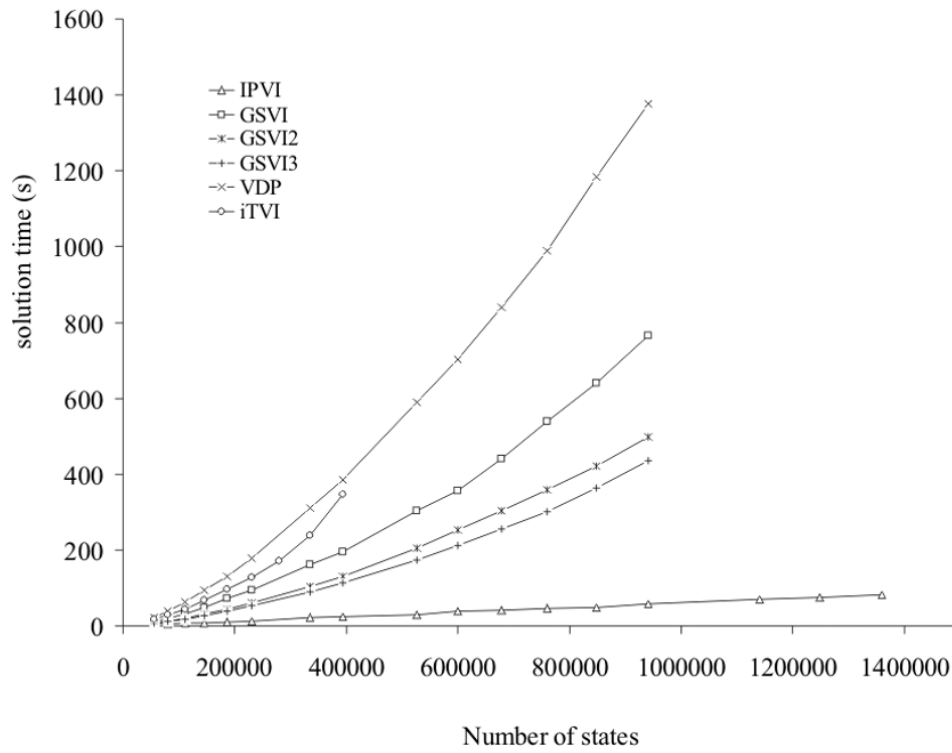


Figure 1 Solution time as a function of the number of states for our algorithm (IPVI), three accelerated variants of the classical VI, VDP (Vanderbei, 2008) and iTVI (Dibangoye, 2008).

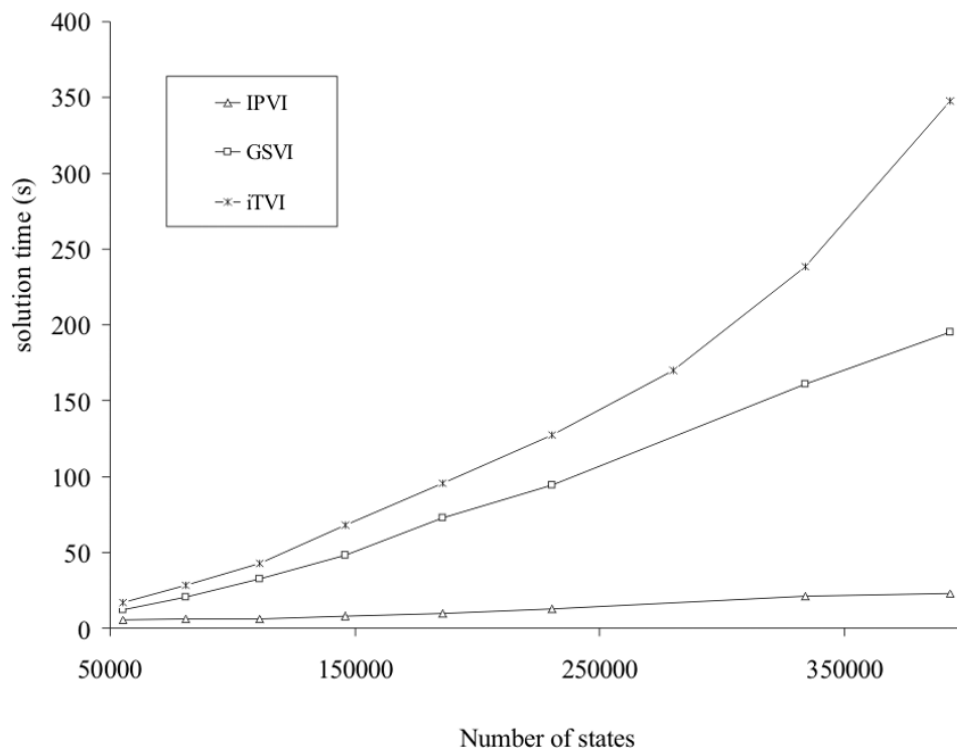


Figure 2 Closer look of the solution time as a function of the number of states for our algorithm (IPVI), GSVI and iTVI (Dibangoye, 2008).

Table 1 Probability of wind direction change. First column indicates old wind direction and first row indicates new wind direction.

	<i>N</i>	<i>NE</i>	<i>E</i>	<i>SE</i>	<i>S</i>	<i>SW</i>	<i>W</i>	<i>NW</i>
<i>N</i>	0.4	0.3						0.3
<i>NE</i>	0.4	0.3	0.3					
<i>E</i>		0.4	0.3	0.3				
<i>SE</i>			0.4	0.3	0.3			
<i>S</i>				0.4	0.2	0.4		
<i>SW</i>					0.3	0.3	0.4	
<i>W</i>						0.3	0.3	0.4
<i>NW</i>	0.4						0.3	0.3