
Efficient Non-Linear Control by Combining Q-learning with Local Linear Controllers

Hajime Kimura*
Tokyo Institute of Technology
gen@fe.dis.titech.ac.jp

Shigenobu Kobayashi
Tokyo Institute of Technology
kobayasi@dis.titech.ac.jp

Abstract

This paper presents a new approach to reinforcement learning (RL) to solve a non-linear control problem efficiently in which state and action spaces are continuous. In real-world applications, an approach combining discrete RL methods with linear controllers is promising since there are many non-linear control problems that can be decomposed into several local linear-control tasks. We provide a hierarchical RL algorithm composed of local linear controllers and Q-learning, which are both very simple. The continuous state-action space is discretized into an array of coarse boxes, and each box has its own local linear controller as an abstract action. The higher-level of the hierarchy is a conventional discrete RL algorithm that chooses the abstract actions. Each linear controller improves the local control policy by using an actor-critic method. The coarse state-space discretization is a quite simple way to cope with the curse of dimensionality, but often gives rise to non-Markovian effects. In our approach, the local linear controllers make up for these undesirable effects. The algorithm was applied to a simulation of a cart-pole swing-up problem, and feasible solutions are found in less time than those of conventional discrete RL methods.

1 Introduction

Many DP-based reinforcement learning (RL) algorithms approximate the value function and give a greedy policy with respect to the learned value

function. Theoretical results guarantee that several DP-based algorithms will find optimal policies (e.g. [Watkins et al. 92], [Jaakkola et al. 93], [Tsitsiklis et al. 97], etc.) and a great deal of effort has been made on the techniques to approximate value functions (e.g., CMAC, Neural-Net [Lin 93]). However, it is too expensive to fit highly accurate value functions, particularly in continuous state-action spaces. To overcome this problem, several techniques are proposed (e.g., Parti-game [Moore 95], Coarse grids +environment models +search techniques [Davies et al. 98], interpolation on a coarse grid [Davies 96]). On the other hand, it is shown that a local gradient-ascent search over stochastic policy-space is possible without explicitly computing value or gradient estimates [Williams 92], [Kimura et al. 95], [Kimura et al. 98]. That is to say, the DP-based RL methods have features of a global search in terms of optimality and computational expense, and the gradient RL methods have a feature of a local search. This paper presents a new model-free approach that bridges a gap between DP-based methods and local gradient-ascent methods.

In real-world applications, an approach combining discrete RL methods with linear controllers is promising since there are many non-linear control problems that can be decomposed into several local linear control tasks. Based on this principle, we provide a new algorithm with the key ideas including:

- A hierarchical RL method composed of Q-learning and local linear controllers, which are both very simple.
- Generalization techniques for continuous state-action spaces: *coarse discretization* for the high-level of the hierarchy, and *local continuous-vector* representation for the low-level linear controllers.
- A policy improvement algorithm for the local linear controllers with imperfect value functions.

The coarse state-space quantization is a quite simple

* Interdisciplinary Graduate School of Science and Engineering, Tokyo Institute of Technology, 4259 Nagatsuta Midori-ku Yokohama 226-8502 JAPAN.

way to cope with the curse of dimensionality. But in many cases, it often gives rise to undesirable non-Markovian effects. In our approach, local linear controllers, which use a continuous-vector for the local state representation, make up for these effects.

Although hierarchical algorithms can often yield sub-optimal solutions, there are significant benefits in learning efficiency, search space, and re-use of knowledge [Singh 92], [Kaelbling 93], [Dietterich98], [Sutton et al. 98], [Parr et al. 98]. In general, the low-level of the hierarchy is composed of learning control modules that represent subtasks or abstract actions. Each subtask is defined in terms of termination conditions, and the low-level module is to find a local optimal policy. The RL techniques for discrete semi-Markov decision processes (SMDPs) would be applied to learning of the high-level module that selects a low-level modules as abstract action. In our approach, the low-level modules correspond to linear controllers.

Actor-critic methods [Barto et al. 83] are often used for learning on the linear controllers [Gullapalli 92], [Doya 96]. Generic actor-critic methods can be classified into a local gradient-ascent method with respect to policy space by using learned value functions. But in our hierarchy, the linear controllers cannot hold accurate value functions on account of practical limitations on the computational resources. Fortunately, it is shown that an actor-critic algorithm using eligibility traces in the actor can improve its stochastic policy even though the estimated value function is inaccurate [Kimura et al. 98]. We take advantage of this method for the local policy improvement.

This paper is structured as follows. First, we review the basic notion of RL and control problems in continuous state dynamic systems. In Section 3 we introduce the generalization techniques for continuous state-action spaces and the hierarchical algorithm combining Q-learning with linear controllers. Section 4 presents behavior of the algorithm through simulations of a cart-pole swing-up task. In Section 5, we summarize the paper and discuss some directions of combining other methods for future work.

2 Problem Formulation

2.1 Markov Decision Processes

At each discrete time t , the agent observes x_t containing information about its current state, select action a_t , and then receives an instantaneous reward r_t resulting from state transition in the environment. In general, the reward and the next state may be random, but their probability distributions are assumed to depend only on x_t and a_t in Markov decision processes (MDPs), in which many reinforcement learning

algorithms are studied. In MDPs, the next state y is chosen according to the transition probability p_{xy}^a , and the reward is given randomly according to the expectation r_x^a . but the agent does not know p_{xy}^a and r_x^a ahead of time. The objective of reinforcement learning is to construct a policy that maximizes the agent's performance. A natural performance measure for infinite horizon tasks is the cumulative discounted reward:

$$V_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k} , \quad (1)$$

where the discount factor, $0 \leq \gamma < 1$ specifies the importance of future rewards, and V_t is the value on time t . In MDPs, the value can be defined as:

$$V^\pi(x) = E_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_k | x_0 = x \right] , \quad (2)$$

where E_π denotes the expectation assuming the agent always uses stationary policy π . The objective in MDPs is to find an optimal policy that maximizes the value of each state x defined by Equation 2. In MDPs, The optimal value of taking action a in state x , denoted $Q^*(x, a)$, satisfies the Bellman equations for all x and a :

$$Q^*(x, a) = r_x^a + \gamma \sum_y p_{xy}^a \max_{a'} Q^*(y, a') \quad (3)$$

2.2 A Non-Linear Control Problem

We are given learning control of non-linear dynamic systems in which :

- State and action spaces are continuous and multidimensional.
- The task could be decomposed into several local linear or bang-bang control tasks.
- The agent does not know dynamics of the environment ahead of time.

Throughout the paper, we use a cart-pole swing-up task for the example.

3 Combining Q-learning with Local Linear Controllers

3.1 Hierarchical Decomposition

In our approach, the agent adopts hierarchical state representation. The higher-level of the hierarchy uses discrete representation of the state variables by coarse quantization. Representative points are situated at the center of the grid's boxes. Each box has its own local linear controllers as discrete abstract actions, that

is, the controllers in a particular box are differently initialized.

Assume that given a n -dimensional state input (x_1, x_2, \dots, x_n) and the corresponding box (B_i) with the representative point at $(b_1^i, b_2^i, \dots, b_n^i)$, then the low-level linear controller gets a local continuous input

$$C = (c_1, c_2 \dots c_n) = (x_1 - b_1^i, x_2 - b_2^i, \dots, x_n - b_n^i). \quad (4)$$

Figure 1 shows an example of the state representation in which the state space is two-dimensional.

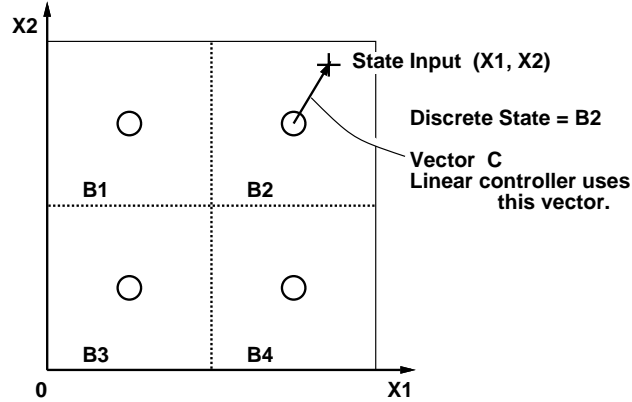


Figure 1: An example of the hierarchical state representation in a 2-dimensional task. Small circles at the center of the grid's boxes denote its representative points.

3.2 Event-driven Decision Making

In dynamic control, an optimal policy often selects the same action (or holds similar continuous action) for a certain period of time. Then, uniform regions are likely to exist in the state space where all of the states have the same (or similar) optimal action. In our approach, the coarse grids approximate the uniform regions. The high-level decision maker selects abstract action each time when the continuous state input is going across the boundary of the box, or a certain period of time passes. This model of abstract action can be seen as an extension of Markov options [Sutton et al. 98] in which the policy is given by linear control and the termination condition is given by the boundary of the box.

Figure 2 helps to define the big-steps that the high-level learner takes. If all the linear controllers hold stationary policy, the high-level learner is to solve a semi-Markov decision problem.

When a linear controller (option) o is taken, then primitive actions are selected according to its linear control rule π^o until the next decision making occurs on the high-level hierarchy. In this paper, we will use the following notation based on [Sutton et al. 98]. Let $\pi^o(x)$

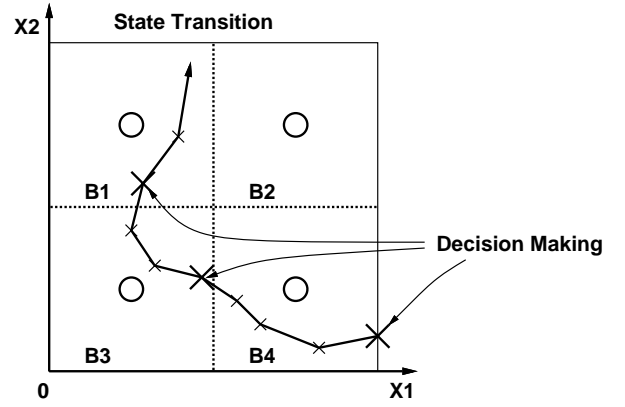


Figure 2: The big-steps on the high-level learner. This modeling is closely related to semi-Markov decision processes (SMDPs). The high-level decisions are only allowed at events which occur whenever the state transition goes across the boundary of the box, or a certain period of time passes.

be the primitive action taken in state x according to o . Assume that current o terminates at next time step according to the probability $T(x, a)$, where $a = \pi^o(x)$. Then, we can define the optimal option-value function recursively by

$$Q^*(x, o) = r_x^a + \gamma(1 - T(x, a)) \left(\sum_y p_{xy}^a Q^*(y, o) \right) + \gamma T(x, a) \left(\sum_y p_{xy}^a \max_{o'} Q^*(y, o') \right), \quad (5)$$

where $a = \pi^o(x)$. If all states x, y are discrete and π^o is stationary, then Equation 5 can be solved by Intra-Option value learning [Sutton et al. 98]. Unfortunately, the state spaces are continuous and it is too expensive to learn the accurate value function. In the next section, we will present a new method that uses Q-values averaged in the coarse regions.

3.3 A Learning Algorithm

Figure 3 gives an overview of the hierarchical algorithm. The high-level decision maker accumulates Q-values using tables with one entry for each box-option pair. In this case, each box-option value (or box value) $Q(B_i, o)$ and $V(B_i)$ represent approximation of averaged values weighted by the visiting frequency over the coarse regions, i.e., $V(B_i) \simeq \sum_{x \in B_i} U^\pi(x) V^\pi(x)$, and $Q(B_i, o) \simeq \sum_{x \in B_i} U^\pi(x) Q^*(x, o)$, where B_i denotes the box, $U^o(x)$ denotes the probability of occupying state x under the policy π . The learner also accumulates state value $V(B_i)$. Unfortunately, this value function does not satisfy Bellman equations in

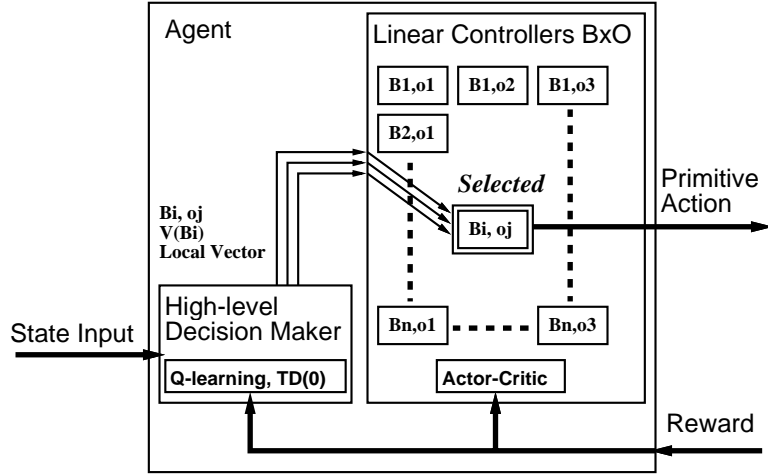


Figure 3: The hierarchical structure of the algorithm. $B_1, B_2, \dots, B_i \dots B_n$ denote the state boxes, and every small rectangle labeled “ B_*, o_* ” denotes a linear controller. In this figure, a particular box B_i has three controllers which are differently initialized.

the strict sense. However, we try to provide a good approximation algorithm.

Assume that the agent selects o_t at time t in the region $B(t)$ and makes next decision at time $t + k + 1$ in the region $B(t + k + 1)$. Then we approximate values at each time by:

$$\begin{aligned}
 E\{V_t^\pi\} &\simeq E\{r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} \dots + \gamma^k r_{t+k}\} \\
 &\quad + \gamma^{k+1} V(B(t+k+1)), \\
 E\{V_{t+1}^\pi\} &\simeq E\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots + \gamma^{k-1} r_{t+k}\} \\
 &\quad + \gamma^k V(B(t+k+1)), \\
 &\quad \vdots \\
 E\{V_{t+k}^\pi\} &\simeq E\{r_{t+k}\} + \gamma V(B(t+k+1)).
 \end{aligned}$$

We approximate the value of the box B_i according to averaging these values:

$$V(B_i) = \frac{1}{k+1} (V_t^\pi + V_{t+1}^\pi \dots + V_{t+k}^\pi). \quad (6)$$

Then, the TD-error for $V(B)$ in the high-level module is given by

$$\begin{aligned}
 \text{TD-error} &= \frac{1}{k+1} \left(\sum_{i=0}^k r_{t+i} \sum_{j=0}^i \gamma^j \right) \\
 &\quad + \frac{1}{k+1} \left(\sum_{i=0}^k \gamma^{i+1} \right) V(B(t+k+1)) - V(B(t)).
 \end{aligned} \quad (7)$$

This result leads us to the learning algorithm for the high-level module.

A Learning Method for High-Level Module:

To estimate box values, the algorithm updates by

$$V(B(t)) \leftarrow V(B(t)) + \alpha_v (\text{TD-error}), \quad (8)$$

where α_v is a learning coefficient, and TD-error is given by (7). Similarly, the box-option value is updated by

$$\begin{aligned}
 Q(B(t), o_t) &\leftarrow Q(B(t), o_t) \\
 &\quad + \alpha_v \left[\frac{1}{k+1} \left(\sum_{i=0}^k r_{t+i} \sum_{j=0}^i \gamma^j \right) \right. \\
 &\quad + \frac{1}{k+1} \left(\sum_{i=0}^k \gamma^{i+1} \right) \max_{o'} Q(B(t+k+1), o') \\
 &\quad \left. - Q(B(t), o_t) \right].
 \end{aligned} \quad (9)$$

Although this algorithm is ad hoc, it works very well for two reasons:

1. Policies that are greedy for Q-functions are optimal if the Q-functions are a certain degree of approximation [Singh et al. 94], [Heger 96].
2. In our approach, learning in the linear controllers makes up for undesirable effects.

A Learning Method for Linear Controllers:

The high-level decision maker selects one of $B \times O$ linear controllers according to some exploration strategy, where B denotes the number of the state boxes, O is the number of possible options per each box. When the state input vector is n -dimensional, the linear controllers have $n + 1$ parameters $W = (w_1, w_2, \dots, w_{n+1})$ per dimension of the output vector.

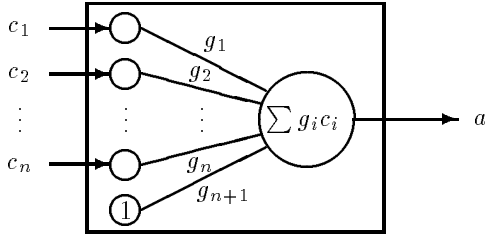


Figure 4: An instance of the linear controller. g_i denotes a feedback gain which is characterized by w_i according to $g_i = w_i + N(0, \sigma^2)$. The controller learns w_i to achieve good control.

The controllers learn W to make good control. A linear controller o determines its state feedback-gain parameters $(g_1, g_2, \dots, g_{n+1})$ as to the control rule π^o by using W . After that, primitive actions are selected according to π^o until the termination of the option. The linear controllers adopt an actor-critic architecture to learn feedback-gain parameters. The actor-critic provides continuous action as the feedback-gain, and makes use of $V(B_i)$ for the critic. The output of the critic is too coarse to calculate value gradient from the learned value function, because the approximation function $V(B_i)$ is flat in the same state region. Therefore we adopt a slightly modified actor-critic algorithm using eligibility traces in the actor [Kimura et al. 98], that can improve its policy even though the estimated value function is inaccurate.

Figure 4 shows an instance of the actor-critic algorithm in which the action is of one dimension. In the beginning of the option, the actor selects its feedback-gain parameters $(g_1, g_2, \dots, g_{n+1})$ according to $g_i = w_i + N(0, \sigma^2)$ for all $i = 1, 2, \dots, n + 1$, where $N(0, \sigma^2)$ denotes the normal distribution. At each time step $t < t + i \leq t + k$, primitive actions are selected until the termination of the option according to

$$a = c_1 g_1 + c_2 g_2 + \dots + c_n g_n + g_{n+1}, \quad (10)$$

where $C = (c_1, c_2, \dots, c_n)$ is the local continuous vector at the time $t + i$ as shown in Equation 4. When the option terminates at time $t + k + 1$ in the region $B(t + k + 1)$, all linear controllers updates according to

$$\begin{aligned} Trace_i &\leftarrow Trace_i + (g_i - w_i) \\ w_i &\leftarrow w_i + \alpha_{ac} (\text{TD-error}) Trace_i \\ Trace_i &\leftarrow \frac{1}{k+1} \left(\sum_{j=0}^k \gamma^{j+1} \right) Trace_i \end{aligned}$$

, where $Trace_i$ denotes an eligibility trace on w_i , and α_{ac} is a coefficient of learning. TD-error is shown in Equation 7. The eligibility of w_i corresponds to $(g_i - w_i)$ that is similar to *Gaussian unit* [Williams 92] and Gullapalli's *neural reinforcement learning unit for continuous action* [Gullapalli 92].

Initial feedback-gain parameters are all set to zero, except that the linear controllers in the same state region have different initial offsets of the primitive actions, i.e., every w_{n+1} are different. Owing to this initializing, the agent can execute various actions even if the linear controllers do not move the parameters at all.

4 Applying to a Cart-Pole Swing-Up Task

The behavior of this algorithm is demonstrated through a computer simulation of a cart-pole swing-up task. We modified the cart-pole problem described in [Barto et al. 83] so that the action is taken to be continuous.

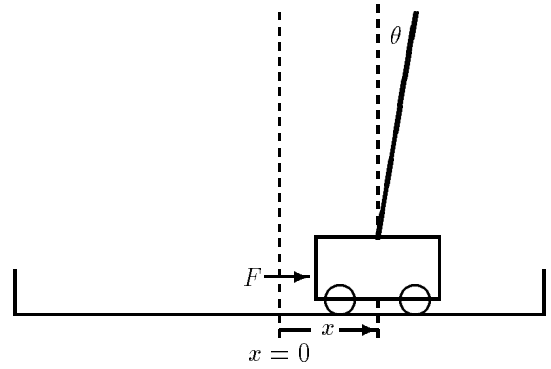


Figure 5: A simulation model of the cart-pole task.

4.1 Details of the Cart-Pole Simulation

The dynamics are modeled by

$$\begin{aligned} \ddot{\theta} &= \frac{g \sin \theta + \cos \theta \left(\frac{-F - m \ell \dot{\theta}^2 \sin \theta + \mu_c s g n(\dot{x})}{M + m} \right) - \frac{\mu_p \dot{\theta}}{m \ell}}{\ell \left(\frac{4}{3} - \frac{m \cos^2 \theta}{M + m} \right)}, \\ \ddot{x} &= \frac{F + m \ell \left(\dot{\theta}^2 \sin \theta - \ddot{\theta} \cos \theta \right) - \mu_c s g n(\dot{x})}{M + m}, \end{aligned}$$

where $M = 1.0(\text{kg})$ denotes mass of the cart, $m = 0.1(\text{kg})$ is mass of the pole, $2\ell = 1(\text{m})$ is the length of the pole, $g = 9.8(\text{m}/\text{sec}^2)$ is the acceleration of gravity, $F(\text{N})$ denotes the force applied to cart's center of mass, $\mu_c = 0.05$ and $\mu_p = 0.01$ are coefficient of friction of the pole and the cart respectively. In this simulation, we use discrete-time system to approximate these equations, where $\Delta t = 0.02\text{sec}$. At each discrete time step, the agent observes $(x, \dot{x}, \theta, \dot{\theta})$, and controls the force F . The agent can execute action in arbitrary range, but the possible action in the cart-pole system is constrained to lie in the range $[-10, 10](\text{N})$. When the agent chooses an action which does not lie in that

range, the action F is saturated. The system begins with $(x, \dot{x}, \theta, \dot{\theta}) = (0, 0, 3.0, 0)$. When the cart collides with the end of the track ($-3.0 \leq x \leq 3.0$), the cart rebounds from the bumper with a coefficient of rebound 0.2. The agent receives a reward (penalty) signal of

- (a) -1 when the pole falls over $\pm 0.8\pi$ (rad),
- (b) -4 when the cart bounces at the end of the track,
- (c) -3 when $\dot{\theta} < -10$ or $10 < \dot{\theta}$ (rad/sec),
- (d) $+1$ when $-0.133\pi < 0.133\pi$ and $-2 < \dot{\theta} < 2$ (rad/sec).

Note that the domain shown here differs from traditional cart-pole tasks in that the environment has no termination condition.

4.2 Implementation

In this experiment, the state space is normalized as $(x, \dot{x}, \theta, \dot{\theta}) = (\pm 3.0 \text{ m}, \pm 10 \text{ m/sec}, \pm \pi \text{ rad}, \pm 10 \text{ rad/sec})$ into $(\pm 0.5, \pm 0.5, \pm 0.5, \pm 0.5)$. The agent discretizes the normalized state space evenly into $3 \times 3 \times 5 \times 5 = 225$ or $3 \times 3 \times 6 \times 6 = 324$ boxes, and attempts to store in each box $V(B_i)$ and $Q(B_i, o)$. Each linear controller accumulates five feedback-gain parameters w_1, w_2, \dots, w_5 , for the continuous-state input is 4-dimensional. On each state box, the high-level decision maker selects one of 3 (or 4) options. When 3 options are within the box, the initial parameters are set to $w_5 = -1/3, w_5 = 0, w_5 = 1/3$ respectively. When 4 options are within the box, the initial parameters are set to $w_5 = -3/8, w_5 = -1/8, w_5 = 1/8, w_5 = 3/8$ respectively. In the linear controllers, an action a is generated by Equation 10, where the normal distribution follows $\sigma = 0.5$. Then the force F is executed according to $F = a \times 20$.

The high-level hierarchy makes decision each time when the continuous state input moves to a different box or the state input keeps within the same box for 2 seconds (100 steps). It uses ϵ -greedy exploration strategy, that is, choosing greedy options with probability 0.99^k , where k is the period of time steps on the previous option. The learning rate α_v is set to 0.3 in $3 \times 3 \times 6 \times 6$ boxes or to 0.1 in $3 \times 3 \times 5 \times 5$ boxes. The learning coefficient of the actor-critic is set to $\alpha_{ac} = 0.01$. The discount rate is set to $\gamma = 0.98$.

4.3 Simulation Results

Figure 6, 9, 12, and 15 show the on-line performance with using different grids and different number of abstract actions. The performance measure is the occupancy rate of the current state in which the system gives positive reward. The rate is calculated by using 20 independent runs. The results show that the proposed methods achieved learning to gain positive rewards. However, Q-learning-only methods in which

linear controllers do not learn (i.e., $\alpha_{ac} = 0$) couldn't learn it at all. It is clear that learning on linear controllers has these effects. The increase of the performance with $3 \times 3 \times 5 \times 5$ boxes is better than with $3 \times 3 \times 6 \times 6$ boxes. One reason for this is an effect of decreasing the number of boxes by the coarse discretization. The other is that the location of the boundaries of the boxes fits for this task fortuitously. Since the agents retained the exploration strategy for all learning steps, the performance could not approach to one. Figure 7, 10, 13 and 16 show examples of trajectories on learned greedy policies after 7,200,000 steps (40 hours on simulation time). The policies were feasible, but no agent could obtain optimal swing-up behavior. The reason is that the growth of swing-up behavior is slowed down for decreasing in opportunity of swing-up, because the agent makes progress in keeping the pole vertical.

Figure 8, 11, 14 and 17 show the corresponding state-action flow respectively. The behavior seems to be a mixture of bang-bang control and linear control. Around $\theta = 0$, we should notice that a bang-bang control rule was found by the algorithms with $3 \times 3 \times 6 \times 6$ boxes whereas a linear control rule was found by the algorithms with $3 \times 3 \times 5 \times 5$ boxes. Anyway, the algorithm finds preferable solutions.

Figure 18 and 19 show the on-line performance of conventional Q-learning using $3 \times 3 \times 10 \times 5$ and $3 \times 3 \times 12 \times 6$ boxes, which are just twice the boxes. The performance is slightly better than that of Q-learning only in Figure 6, 9, 6 and 9, but they could not learn the control to keep the pole standing. The results show that the task is not trivial.

5 Discussion

Evaluation: We cannot conclude superiority of the proposed method, because the complexity of our approach is obviously larger than that of conventional Q-learning in the experiment. The experiment shows that the learning of the local linear controllers can make up for lack of control ability in the high-level decision maker.

Effects of the Partitioning Location: In our test cases, different discretizations led the algorithms to learning quite different control rules, especially around $\theta = 0$. These things make it clear that the shape and location of the quantized regions have great effects on the performance. Joint use of variable-grid methods (e.g. parti-game [Moore 95]), would be needed for generating discrete representation rather than the use of the fixed grid.

Time Intervals of High-level Decision Making: In our approach, time intervals of decision making on

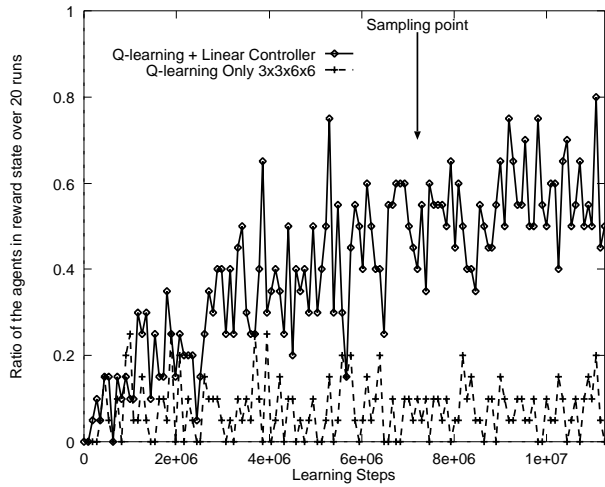


Figure 6: On-line performance using 3x3x6x6 grids and 4 abstract actions.

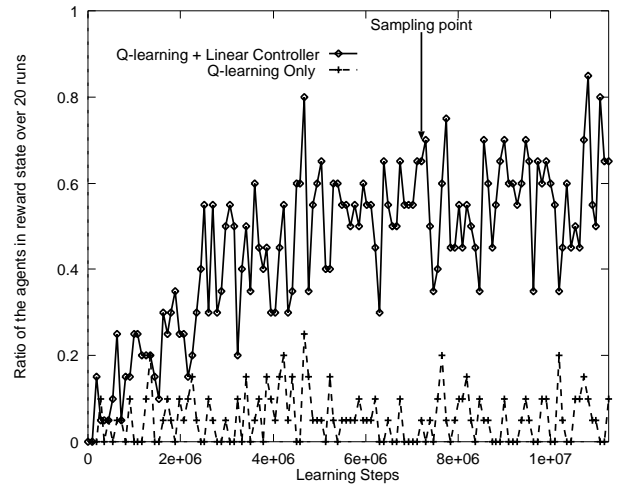


Figure 9: On-line performance using 3x3x6x6 grids and 3 abstract actions.

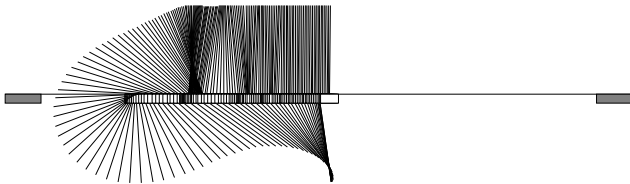


Figure 7: An example of swing-up behavior on 3x3x6x6 grids and 4 actions.

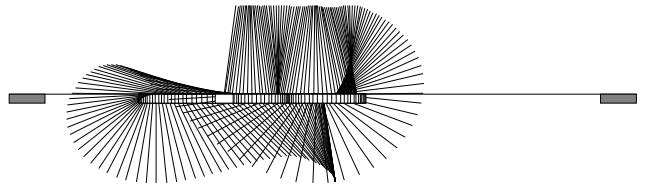


Figure 10: An example of swing-up behavior on 3x3x6x6 grids and 3 actions.

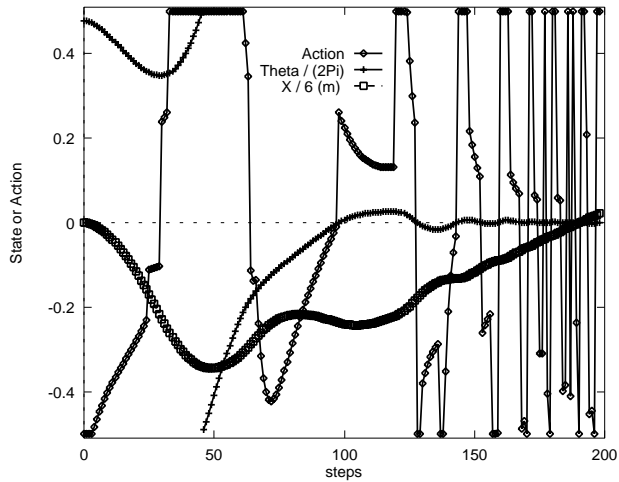


Figure 8: Detail of the behavior.

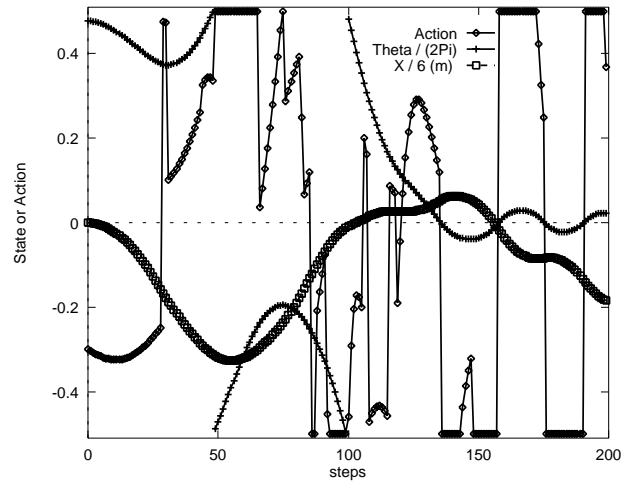


Figure 11: Detail of the behavior.

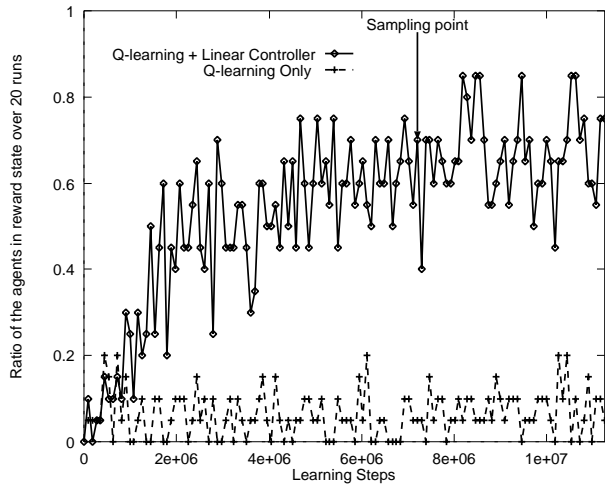


Figure 12: On-line performance using 3x3x5x5 grids and 4 abstract actions.

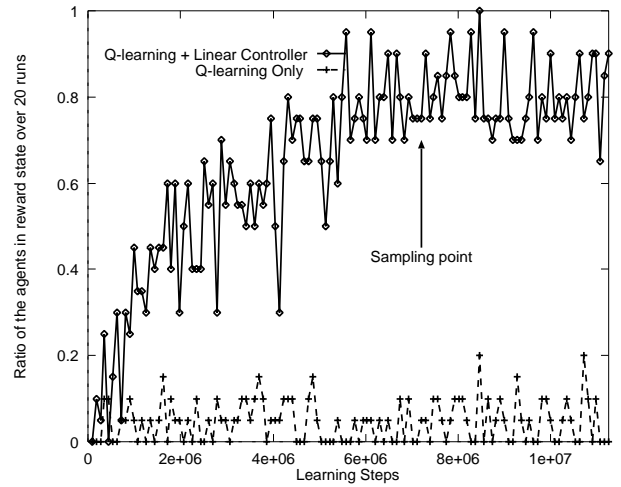


Figure 15: On-line performance using 3x3x5x5 grids and 3 abstract actions.

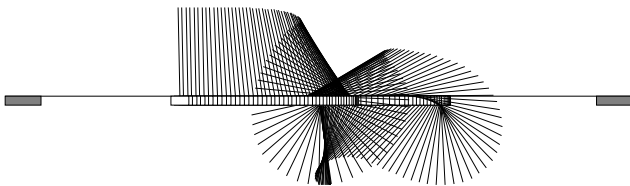


Figure 13: An example of swing-up behavior on 3x3x5x5 grids and 4 actions.

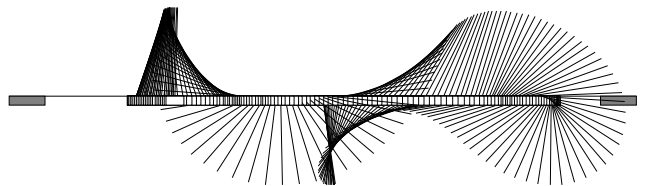


Figure 16: An example of swing-up behavior on 3x3x5x5 grids and 3 actions.

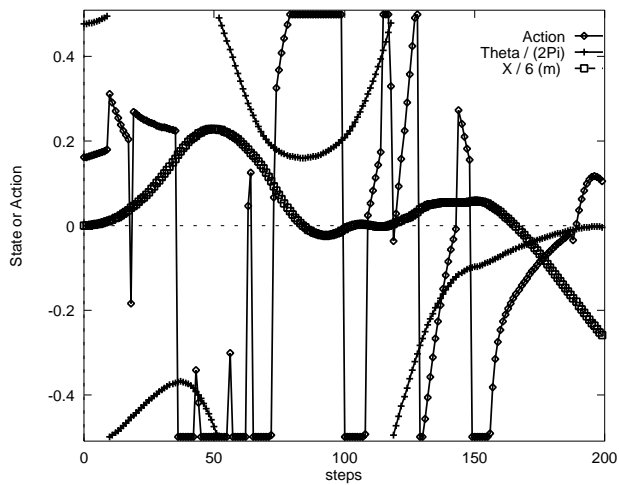


Figure 14: Detail of the behavior.

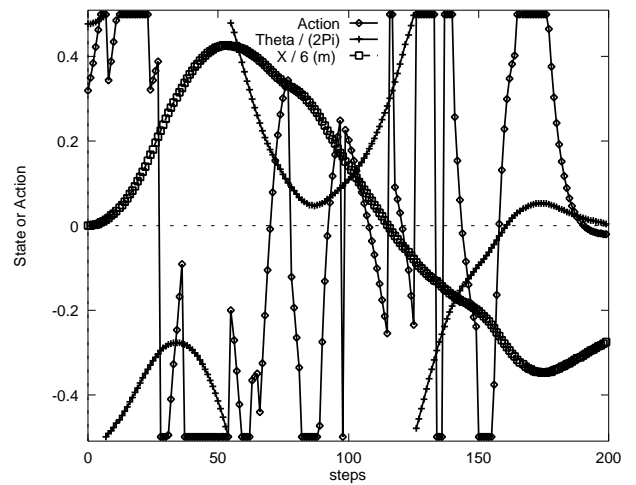


Figure 17: Detail of the behavior.

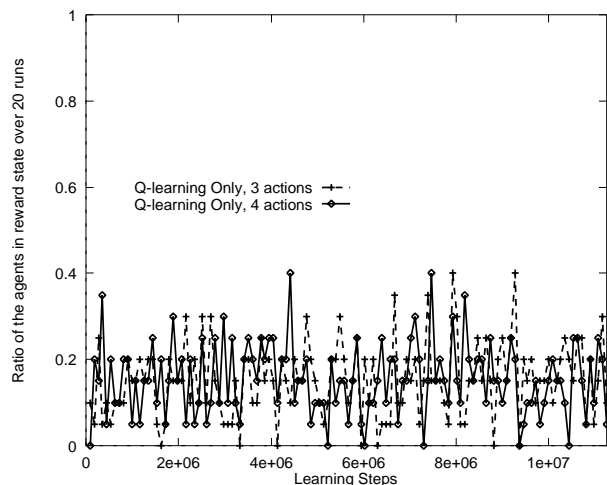


Figure 18: On-line performance of conventional Q-learning using $3 \times 3 \times 10 \times 5$ grids and 3 (or 4) abstract actions. It is just twice as many as the boxes in Figure 12 and 15.

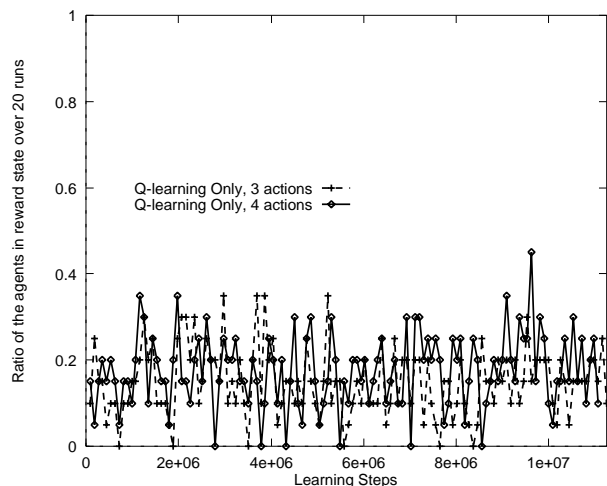


Figure 19: On-line performance of conventional Q-learning using $3 \times 3 \times 12 \times 6$ grids and 3 (or 4) abstract actions. It is just twice as many as the boxes in Figure 6 and 9.

the high-level hierarchy are mostly owing to the size of the state boxes. It can be seen as an adaptive choice of time intervals in continuous-time domains [Pareigis 98] [Buckland et al. 93].

Learning with Extremely Coarse Value Approximation: Many DP-based (RL) algorithms are motivated by a desire for finding highly accurate value functions. However, it costs too much memory to approximate such functions in many cases. In contrast, our approach does not stick to finding accurate value functions, because the policy can be improved by a gradient-ascent search without explicitly computing value estimates. However, shown in the experiments, it costs many experiences instead of the memory required in DP-based methods.

Combining with Model-based Methods: Since the proposed algorithm in this paper is a model-free and memory-less method, many model-based methods can make use of it easily. It is interesting to note that [Davies et al. 98] proposed a model-based approach that uses environment models to make up for the poor ability of the value function approximation.

6 Conclusions

This paper presented a hierarchical RL algorithm composed of Q-learning and local linear controllers to solve a non-linear control problem in which state and action spaces are continuous. It is a hybrid method between DP-based value estimation and policy improvement by gradient-ascent without value estimation. The continuous state-action space is discretized into an array of coarse boxes, and roughly the high-level hierarchy estimates the value functions over the discrete space to find a globally preferable policy. The local linear controllers are to improve the policy by stochastic gradient-ascent. This method does not need to fit accurate value functions, therefore it may be promising to overcome the curse of dimensionality. The algorithm was applied to a simulation of a cart-pole swing-up problem, and better solutions are found than those of traditional discrete RL methods.

Acknowledgements

We would like to thank reviewers for many helpful comments and suggestions.

References

- [Barto et al. 83] Barto, A. G., Sutton, R. S. and Anderson, C. W.: Neuronlike Adaptive Elements That Can Solve Difficult Learning Control Problems, *IEEE Transactions on Systems, Man, and Cybernetics*, vol. SMC-13, no.5, September/October 1983, pp. 834-846.

- [Buckland et al. 93] Buckland, K. M. and Lawrence, P. D.: Transition Point Dynamic Programming, *Advances in Neural Information Processing Systems 6*, pp. 639–646 (1993).
- [Davies 96] Davies, S.: Multidimensional Triangulation and Interpolation for Reinforcement Learning, *Advances in Neural Information Processing Systems 9*, pp. 1005-1011 (1996).
- [Davies et al. 98] Davies, S., Ng, A. Y. & Moore, A.: Applying Online Search Techniques to Continuous-State Reinforcement Learning, *15th National Conference on Artificial Intelligence*, pp. 753–760 (1998).
- [Dietterich98] Dietterich, T. G.: The MAXQ Method for Hierarchical Reinforcement Learning, *Proceedings of the 15th International Conference on Machine Learning*, pp. 118–126 (1998).
- [Doya 96] Doya, K. : Efficient Nonlinear Control with Actor-Tutor Architecture, *Advances in Neural Information Processing Systems 9*, pp. 1012–1018 (1996).
- [Gullapalli 92] Gullapalli, V.: Reinforcement Learning and Its Application to Control, *PhD Thesis*, University of Massachusetts, Amherst, COINS Technical Report 92-10 (1992).
- [Heger 96] Heger, M.: The Loss from Imperfect Value Functions in Expectation-Based and Minimax-Based Tasks, *Machine Learning, 22*, pp. 197–225 (1996).
- [Jaakkola et al. 93] Jaakkola, T., Jordan, M. I. & Singh, S. P.: On the Convergence of Stochastic Iterative Dynamic Programming Algorithms, *Neural Computation 6*, pp.1185–1201 (1993).
- [Kaelbling 93] Kaelbling, L. P.: Hierarchical Learning in Stochastic Domains: Preliminary Results, *Proceedings of the 10th International Conference on Machine Learning*, pp. 167–173 (1993).
- [Kaelbling et al.96] Kaelbling, L. P., & Littman, M. L., & Moore, A. W.: Reinforcement Learning: A Survey, *Journal of Artificial Intelligence Research*, Vol. 4, pp. 237–277 (1996).
- [Kimura et al. 95] Kimura, H., Yamamura, M. and Kobayashi, S.: Reinforcement Learning by Stochastic Hill Climbing on Discounted Reward, *Proceedings of the 12th International Conference on Machine Learning*, pp. 295–303 (1995).
- [Kimura et al. 98] Kimura, H. & Kobayashi, S.: An Analysis of Actor/Critic Algorithms using Eligibility Traces: Reinforcement Learning with Imperfect Value Function, *15th International Conference on Machine Learning*, pp.278–286 (1998).
- [Lin 93] Lin, L. J.: Scaling Up Reinforcement Learning for Robot Control, *Proceedings of the 10th International Conference on Machine Learning*, pp. 182–189 (1993).
- [Lin et al. 96] Lin, C. J. and Lin, C. T.: Reinforcement Learning for An ART-Based Fuzzy Adaptive Learning Control Network, *IEEE Transactions on Neural Networks*, Vol.7, No. 3, pp. 709-731 (1996).
- [Moore 95] Moore A. W. & Atkeson, C. G.: The Partigame Algorithm for Variable Resolution Reinforcement Learning in Multidimensional State-spaces, *Machine Learning 21*, pp. 199-233 (1995).
- [Pareigis 98] Pareigis, S.: Adaptive choice of grid and time in reinforcement learning, *Advances in Neural Information Processing Systems 10*, pp. 1036–1042 (1998).
- [Parr et al. 98] Parr, R. & Russell, S.: Reinforcement Learning with Hierarchies of Machines, *Advances in Neural Information Processing Systems 10*, pp. 1043–1049 (1998).
- [Singh 92] Singh, S. P.: Transfer of Learning by Composing Solutions of Elemental Sequential Tasks, *Machine Learning 8*, pp. 323-339 (1992).
- [Singh et al. 94] Singh, S. P. and Yee, R. C.: An Upper Bound on the Loss from Approximate Optimal-Value Functions, *Machine Learning, 16*, pp. 227-233 (1994).
- [Singh 96] Singh, S. P., & Sutton, R.S.: Reinforcement Learning with Replacing Eligibility Traces, *Machine Learning 22*, pp. 123-158 (1996).
- [Sutton 88] Sutton, R. S.: Learning to Predict by the Methods of Temporal Differences, *Machine Learning 3*, pp. 9-44 (1988).
- [Sutton 90] Sutton, R. S.: Reinforcement Learning Architectures for Animats, *Proceedings of the 1st International Conference on Simulation of Adaptive Behavior*, pp. 288-295 (1990).
- [Sutton et al. 98] Sutton, R. S. & Barto, A.: Reinforcement Learning: An Introduction, *A Bradford Book*, The MIT Press (1998).
- [Sutton et al. 98] Sutton, R. S., Precup, D. & Singh, S.: Intra-Option Learning about Temporary Abstract Actions, *Proceedings of the 15th International Conference on Machine Learning*, pp. 556–564 (1998).
- [Tsitsiklis et al. 97] Tsitsiklis, J. N., & Roy, B. V.: An Analysis of Temporal-Difference Learning with Function Approximation, *IEEE Transactions on Automatic Control*, Vol.42, No.5, pp. 674–690 (1997).
- [Watkins et.al 92] Watkins, C. J. C. H., & Dayan, P.: Technical Note: *Q*-Learning, *Machine Learning 8*, pp. 55-68 (1992).
- [Williams 92] Williams, R. J.: Simple Statistical Gradient Following Algorithms for Connectionist Reinforcement Learning, *Machine Learning 8*, pp. 229-256 (1992).