

Constrained reinforcement learning from intrinsic and extrinsic rewards

Eiji Uchibe

Okinawa Institute of Science and Technology
12-22 Suzaki, Uruma, Okinawa, Japan
uchibe@oist.jp

Kenji Doya

ATR Computational Neuroscience Labs
2-2-2 Hikaridai, "Keihanna Science City," Kyoto 619-0288, Japan
doya@oist.jp

Abstract—The main objective of a standard reinforcement learner is usually defined as maximization of a scalar reward function given externally from the environment. On the other hand, an intrinsically motivated reinforcement learner creates an intrinsic reward function from its own criteria such as curiosity, prediction error, and learning progress. This paper proposes a novel approach to deal with both intrinsic and extrinsic rewards for reinforcement learning from a viewpoint of constrained optimization problem. The extrinsic rewards construct inequality constraints to the stochastic policy while the intrinsic reward determines the current objective function for the learning agent. By integrating policy gradient reinforcement learning algorithms and techniques used in nonlinear programming, our proposed method, named the Constrained Policy Gradient Reinforcement Learning (CPGRL), maximizes the long-term average intrinsic reward under the inequality constraints induced by the extrinsic rewards. The CPGRL is successfully applied to a simple MDP problem and a control task of a robot arm.

Index Terms—Intrinsic and extrinsic rewards, reinforcement learning.

I. INTRODUCTION

The main objective of the learning agent is usually determined by experimenters. In the case of reinforcement learning [1], an appropriate reward function should be designed for each task through a trial-and-error process. It is still important to implement learning algorithms that can efficiently improve the learning capabilities, but the principles for designing the appropriate reward functions become important more and more in the future. In many cases, extrinsic rewards are zero everywhere except for a few important points that correspond to the important events. Although designing such a sparse reward function is easier than designing a dense one, the sparse rewards prevent the learning agent to learn efficiently. On the contrary, the intrinsic reward is regarded as a dense reward functions which gives non-zero rewards most of the time because it is usually computed from the agent's internal information such as sensory inputs. Although the intrinsic reward is generally task-independent, it plays an important role for designing an open-ended system.

Recently, learning algorithms with intrinsic rewards have been studied by several researchers. Barto and his colleagues

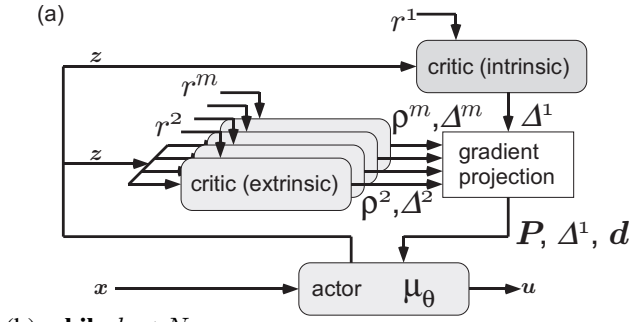
[2], [3], [4] proposed an algorithm for intrinsically motivated reinforcement learning based on the theory of *options* [5]. Meeden *et al.* [6] realized that a simulated robot tracked a moving decoy robot with the rewards based on the error of its own prediction. Oudeyer and Kaplan [7] adopted progress of prediction learning as intrinsic rewards and showed that behavior evolution of the Sony's four-legged robot, AIBO, were realized by step-by-step learning. However, most previous studies did not discuss the negative effects of extrinsic rewards on intrinsically motivated learning suggested by Deci and Flaste [8]. It is still unclear how extrinsic rewards can help or hinder the learning process.

As the first step towards this problem, this paper deals with the interaction between intrinsic and extrinsic rewards from a viewpoint of constrained optimization problems. The learning agent tries to maximize the long-term average intrinsic reward under the inequality constraints given by extrinsic rewards. We propose a new framework termed the *Constrained Policy Gradient Reinforcement Learning* (CPGRL) consisting of a Policy Gradient Reinforcement Learning (PGRL) algorithm [9], [10], [11] and a gradient projection method [12]. Since The PGRL algorithms can estimate the gradients of the expected average rewards with respect to the policy parameters, they are nicely integrated with the gradient projection method. Although constrained Markov Decision Process (MDP) problems are previously studied based on linear programming techniques [13], [14], their methods do not suit our case because the state transition probabilities are known, and because it can not be easily extended to continuous state and action spaces. In order to evaluate the CPGRL we conduct two simulations: a simple MDP problem with three states and a control task of a robotic arm.

II. CPGRL: CONSTRAINED POLICY GRADIENT REINFORCEMENT LEARNING FROM INTRINSIC AND EXTRINSIC REWARDS

A. Formulation

At each time step, an agent observes a state $\mathbf{x} \in \mathbb{X}$ and executes an action $\mathbf{u} \in \mathbb{U}$ with probability $\mu_{\theta}(\mathbf{x}, \mathbf{u})$, where $\mu_{\theta} : \mathbb{X} \times \mathbb{U} \rightarrow [0, 1]$ is a stochastic policy parameterized by



- (b) **while** $k < N_K$
- 1) set $z_0 = \mathbf{0}$ and $\Delta^i = \mathbf{0}$ for all i .
 - 2) **while** $t < N_T$
 - a) observe x_t and execute u_t .
 - b) receive the rewards r_t
 - c) estimate the average rewards and their gradients.
 - 3) store the estimated average rewards
 - 4) update the policy parameter.

Fig. 1. Actor-critic architecture for learning from intrinsic and extrinsic rewards. (a) Block diagram. (b) Algorithm. N_K and N_T denote the number of episodes and steps, respectively.

the n -dimensional vector $\theta \in \mathbb{R}^n$. The agent calculates an intrinsic reward r_t^1 and extrinsic rewards r_t^i ($i = 2, 3, \dots, m$) at time t , which depend on the state and the action. Let $r_t^i = r^i(x_t, u_t)$ and $r_t = [r_t^1 \ r_t^2 \ \dots \ r_t^m]^\top$ denote respectively the immediate reward at time t and the vectorized representation. The operation $^\top$ means the vector/matrix transpose.

The objective for the agent is to find the policy parameter θ that maximizes an average reward

$$g^1(\theta) = \lim_{T \rightarrow \infty} \mathbb{E}_\theta \left[\frac{1}{T} \sum_{t=1}^T r_t^1 \right] \quad (1)$$

under the constraints determined by the extrinsic rewards

$$g^i(\theta) = \lim_{T \rightarrow \infty} \mathbb{E}_\theta \left[\frac{1}{T} \sum_{t=1}^T r_t^i \right] \geq G^i, \quad i = 2, \dots, m, \quad (2)$$

where G^i is a threshold for controlling a level of the constraint. It is noted that the inequality constraints (2) are also the functions of the average rewards.

Fig. 1 illustrates the CPGRL system based on the actor-critic architecture [1]. It consists of one actor, multiple critics, and a gradient projection module that computes a projection onto a feasible region, which is the set of points satisfying all the inequality constraints (2). Based on the immediate reward r^i , each critic outputs ρ^i , an estimate of the long-term average reward g^i , and Δ^i , its gradient with respect to the policy parameters. Actor selects the action u according to the stochastic policy $\mu_\theta(x, u)$.

B. Gradient estimates by policy gradient reinforcement learning

The PGRL algorithms have recently been re-evaluated since they are well-behaved with function approximation. There exist several methods to compute the gradient of the average reward Δ^i . In the current implementation, we choose the GPOMDP algorithm [9] and the actor-critic method [10]. Here, we briefly introduce the GPOMDP algorithm when the reward depends on the action as well as the state. According to the current state and action, the function ψ_t is defined by

$$\psi_t(x_t, u_t) \triangleq \frac{1}{\mu_\theta(x_t, u_t)} \frac{\partial \mu_\theta(x_t, u_t)}{\partial \theta}.$$

The agent interacts with the environment, producing a state, action, reward sequence. After receiving experiences $(x_t, u_t, x_{t+1}, u_{t+1}, r_{t+1})$, the GPOMDP updates an eligibility traces $z_t \in \mathbb{R}^n$

$$z_{t+1} = \beta z_t + \psi_t(x_t, u_t)$$

where $\beta \in [0, 1)$ is a discount factor that controls the variance of the gradient estimate. Since z_t is independent of the reward functions, z_t can be used for estimating gradients of different average rewards. Then, all the gradients are updated in the same manner. That is, the gradient of the long-term average reward is approximated by

$$\Delta_{t+1}^i = \Delta_t^i + \frac{1}{t+1} [r_{t+1}^i (z_{t+1} + \psi_{t+1}(x_{t+1}, u_{t+1})) - \Delta_t^i] \quad (3)$$

for all $i = 1, \dots, m$. The estimate of the average reward r^i is updated by

$$\rho_{t+1}^i = \rho_t^i + \alpha_r (r_{t+1}^i - \rho_t^i), \quad (4)$$

where α_r is a positive step-size meta-parameter. Although the GPOMDP can estimate the gradient with less number of parameters, it has a large variance as $\beta \rightarrow 1$.

We also use a simplified method based on the actor critic method [10] that exploits a value function. The gradient of the long-term average reward is calculated by

$$\Delta_{t+1}^i = \Delta_t^i + \frac{1}{t+1} [Q^i(x_t, u_t) \psi(x_t, u_t) - \Delta_t^i], \quad (5a)$$

$$Q^i(x, u) = (w^i)^\top \psi(x, u). \quad (5b)$$

where $Q^i(x, u)$ and w^i denote an approximated state-action value function and a parameter vector, respectively. In order to train w^i , the standard temporal difference method is carried out

$$w_{t+1}^i = w_t^i + \alpha_r \delta_t^i z_{t+1},$$

where the temporal difference δ_t^i is defined by

$$\delta_t^i = r_{t+1}^i - \rho_{t+1}^i + (w_t^i)^\top [\psi_{t+1}(x_{t+1}, u_{t+1}) - \psi_t(x_t, u_t)].$$

Konda's actor-critic requires an additional learning mechanism to approximate $Q(x, u)$, but it can utilize the Markov property.

C. Gradient projection

As described in section II-B, the average rewards and their gradients are obtained at the end of each episode. Next, we apply a gradient projection method to solve the maximization problem with inequality constraints. When k -th episode ends, the policy parameters are update as follows:

$$\theta_{k+1} = \theta_k + \alpha_1 \mathbf{P} \Delta^1 - \alpha_e \mathbf{d} \quad (6)$$

where $\alpha_1, \alpha_e \in [0, 1)$ are learning rates, \mathbf{P} is a matrix that projects Δ^1 into the subspace tangent to the active constraints, and \mathbf{d} is a restoration move for the violating constraints.

In order to estimate \mathbf{P} and \mathbf{d} , a set of indices of the active inequality constraints is defined by

$$\mathcal{A} = \{i \mid \rho^i - G^i \leq 0, i = 2, \dots, m\},$$

and let $a = |\mathcal{A}|$ denote the number of active constraints. \mathcal{A} is called an active set. If no constraints are active (the case $a = 0$), the solution lies at the interior of the feasible region. In this case, \mathbf{P} and \mathbf{d} are set to the identity matrix and zero vector, respectively. Hereafter, the case $a \neq 0$ is considered. With the outputs from the multiple critics, we define

$$\mathbf{g}_{\mathcal{A}} \triangleq [\rho^{i_1} - G^{i_1} \quad \dots \quad \rho^{i_a} - G^{i_a}]^{\top},$$

$$\mathbf{N}_{\mathcal{A}} \triangleq [\Delta^{i_1} \quad \dots \quad \Delta^{i_a}],$$

where i_a is an index to count the element in \mathcal{A} . The projection matrix and restoration move are given by

$$\mathbf{P} = \mathbf{I} - \mathbf{N}_{\mathcal{A}} (\mathbf{N}_{\mathcal{A}}^{\top} \mathbf{N}_{\mathcal{A}})^{-1} \mathbf{N}_{\mathcal{A}}^{\top}, \quad (7)$$

$$\mathbf{d} = \mathbf{N}_{\mathcal{A}} (\mathbf{N}_{\mathcal{A}}^{\top} \mathbf{N}_{\mathcal{A}})^{-1} \mathbf{g}_{\mathcal{A}}. \quad (8)$$

It should be noted that $\mathbf{P} \mathbf{d} = \mathbf{0}$.

Here, we should note two points when (7) and (8) are computed in the program. At first, $\mathbf{N}_{\mathcal{A}}^{\top} \mathbf{N}_{\mathcal{A}}$ in (7) and (8) is not invertible if the set of active constraint gradients $\{\Delta^{i_j} \mid j = 1, \dots, a\}$ is linearly dependent. In practice, rank deficiency of $\mathbf{N}_{\mathcal{A}}^{\top} \mathbf{N}_{\mathcal{A}}$ is sometimes detected due to the accuracy of numerical computation and/or biased samples. The pseudo-inverse of $\mathbf{N}_{\mathcal{A}}^{\top} \mathbf{N}_{\mathcal{A}}$ should be used if it is not full-rank. In addition we must consider the situation where $\mathbf{P} \Delta^1 = \mathbf{0}$ because it may be possible to modify the parameters. This situation can be detected by using Lagrange multipliers

$$\lambda = (\mathbf{N}_{\mathcal{A}}^{\top} \mathbf{N}_{\mathcal{A}})^{-1} \mathbf{N}_{\mathcal{A}}^{\top} \Delta^1. \quad (9)$$

If λ has no negative components, we have a solution and terminate. Otherwise, the active set is re-evaluated by $\mathcal{A} \leftarrow \mathcal{A} \setminus \{r\}$ where $r = \operatorname{argmax}_{i \in \mathcal{A}} \lambda_i$. After deleting one constraint from the active set, \mathbf{P} and \mathbf{d} are calculated again. See appendix about (7) and (9).

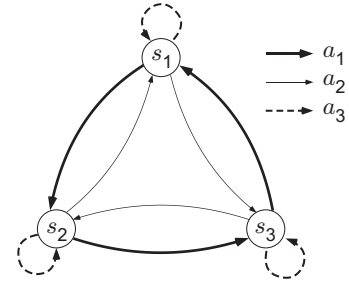


Fig. 2. Simple MDP with three states, three actions, and four reward functions.

III. PERFORMANCE IN A SIMPLE MDP PROBLEM

A. Simulation setup

In order to evaluate the performance of the CPGRL from a viewpoint of constrained optimization problems, we apply the CPGRL to the simple three-state MDP problem shown in Fig. 2. The sets of states and actions are $\{s_1, s_2, s_3\}$ and $\{a_1, a_2, a_3\}$, respectively¹. Each action achieves the intended effect with probability 0.8, but it makes a random transition otherwise. For example, from the state s_1 the action a_1 moves the agent to s_2, s_3, s_1 with probabilities 0.8, 0.1, 0.1, respectively.

One intrinsic reward r^1 and three extrinsic rewards r^2, r^3, r^4 are prepared in this problem;

$$r^1 = \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & 2 \\ 1 & 0 & -1 \end{bmatrix}, \quad r^2 = \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix},$$

$$r^3 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -1 \\ 0 & 0 & 1 \end{bmatrix}, \quad r^4 = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -1 \end{bmatrix},$$

where $r^i = (r_{jk}^i)$ is a reward value of r^i when the action a_k is selected at the state s_j . For instance, the reward vector is $\mathbf{r} = [2 \ 1 \ -1 \ 0]^{\top}$ when the agent selects the action a_3 at the state s_2 . Under these settings, the optimal policy is to select a_1 in each state, and the long-term average reward vector is $[1 \ 0 \ 0 \ 0]^{\top}$. It should be noted that the extrinsic rewards are competitive with each other. As the stochastic policy, we use a lookup table with softmax distribution

$$\mu_{\theta}(x_t, u_t) = \frac{\exp(\theta_{x_t, u_t})}{\sum_{u'} \exp(\theta_{x_t, u'})},$$

yielding a total of nine policy parameters. That is, the policy parameters are assigned such that $\theta_1 = \theta_{s_1, a_1}$, $\theta_2 = \theta_{s_1, a_2}$.

Equation (3) is adopted for estimating the policy gradient. Each policy parameter is randomly initialized with a uniform distribution over the interval $[0, 1]$. The agent uses inequality

¹Let s and a denote the original state and action in this MDP problem while the variables x and u represent the state and action at each time step. Therefore, $x_t \in \{s_1, s_2, s_3\}$ and $u_t \in \{a_1, a_2, a_3\}$.

constraints with $G^2 = G^3 = G^4 = 0$. Other meta-parameters are set as follows: $\alpha_\rho = 0.02$, $\beta = 0.99$, $\alpha_1 = \alpha_e = 0.02$. These values are determined by trial and error. In order to compare the performance, we consider two different approaches named the CONstraints-Based (CONB) and the SUM method. The CONB switches the policy gradient of each reward according to the following condition:

$$\Delta = \begin{cases} \Delta^1 & \text{if } g^i - G^i > 0 \text{ for } i = 2, \dots, m \\ \Delta^j & \text{otherwise, } j = \arg \min_{i \in \mathcal{A}} (\rho^i - G^i). \end{cases}$$

to maximize the average reward of r^1 and to satisfy the constraints. The SUM learns to maximize the average reward of the summation of all rewards

$$r^{\text{sum}} = r^1 + r^2 + r^3 + r^4$$

based on the standard policy gradient method. However, it is expected that the learned parameters does not satisfy the constraints since the SUM does not consider the constraints at all. The agent starts at the state $x_0 = s_1$. The number of episodes and steps are $N_T = 100$ and $N_K = 10000$, respectively. We perform 20 simulation runs.

B. Experimental results

Figs. 3 show the means and the standard deviations of 20 simulation runs obtained by the CPGRL, CONB, and SUM, respectively. The CPGRL found the parameters that satisfy the inequality constraints at the very early stage of learning, and the standard deviations of the average rewards of constraints were very small after 1×10^3 -th episode. The CONB also obtained the policy parameters satisfying constraints, but it took a longer time than the CPGRL. The long-term average reward of r^1 were gradually increased by the CPGRL. Interestingly, the CONB failed to maximize the average reward of r^1 . In addition, we found that the standard deviation estimated by the CONB was larger than that of the CPGRL. Although the SUM obtained the best average reward on r^1 , it failed to find the policy parameters satisfying the constraint on r^3 .

Then, we checked the stochastic policy during the learning process. Fig. 4 shows the evolution of the policy parameters of the CPGRL, CONB, and SUM, respectively. In this figure, all policies lie in the lower-left triangle, and gray-color represents the average reward of r^1 . Although all methods could obtain appropriate action at the state s_1 ($\Pr(a_1 | s_1) = 1$ is optimal), the CONB and SUM obtained inappropriate actions at the states s_2 and s_3 . For example, the SUM leaned to select a_3 at s_2 because the large positive reward ($2 + 1 - 1 + 0 = 2$) was received in this case. Obviously, this violated the constraints on r^3 . The CONB failed to obtain the appropriate action at s_3 . It should be noted that both of a_1 and a_2 did not generate negative rewards in this state. However, the CONB failed to improve the average reward of r^1 because the gradient of r^1 was rarely selected. Since

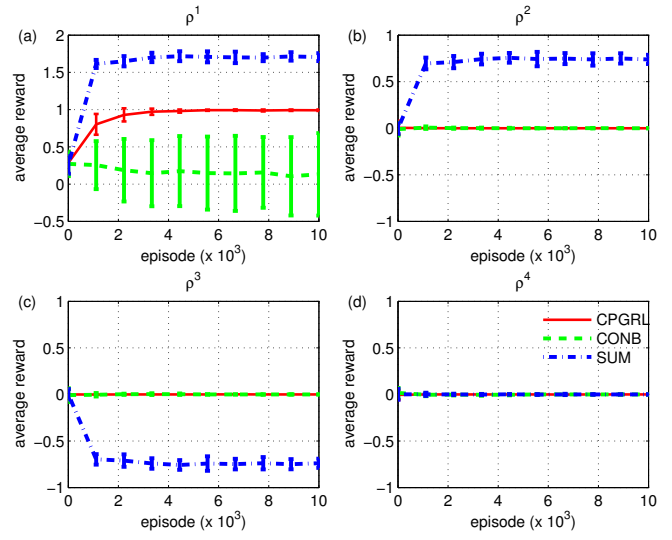


Fig. 3. Transition of the estimated average rewards. (a) ρ^1 , (b) ρ^2 , (c) ρ^3 , and (d) ρ^4 , respectively. These figures show the means and standard deviations of 20 independent runs.

the estimated average reward by (4) is not deterministic, constraints were suddenly violated. The CPGRL successfully obtained the optimal policy that satisfied all constraints in this simulation.

IV. CONTROL TASK OF A ROBOTIC ARM

A. Simulation setup

Then we conduct a control task of a robotic arm to investigate how the thresholds G^i affect the learning processes. Fig. 5 (a) shows an simulated environment. There exists a typical two-link arm that can interact four objects (circle, star, square, and triangle). These four objects are fixed in the environment. The intrinsic reward r^1 is computed from the distance between the position of the end-effector of the arm and the nearest objects. Fig. 5 (b) shows a distribution of r^1 . This dense reward function enables the robotic arm to learn touching behaviors actively. Then, two extrinsic rewards r^2 and r^3 are introduced in this task. The first extrinsic reward gives upper and lower bounds on the joint angles ϕ_1 and ϕ_2 while the second extrinsic reward depends on whether the touched object is appetitive or aversive:

$$r^2 = \begin{cases} 0 & \frac{\pi}{10} \leq \phi_1 \leq \frac{\pi}{2}, \frac{\pi}{4} \leq \phi_2 \leq \frac{3\pi}{4}, \\ -1 & \text{otherwise,} \end{cases}$$

$$r^3 = \begin{cases} 1 & \text{if the object is appetitive,} \\ -1 & \text{if the object is aversive,} \\ 0 & \text{otherwise.} \end{cases}$$

It should be noted that zero reward is given when the robotic arm touches the circular and triangular objects.

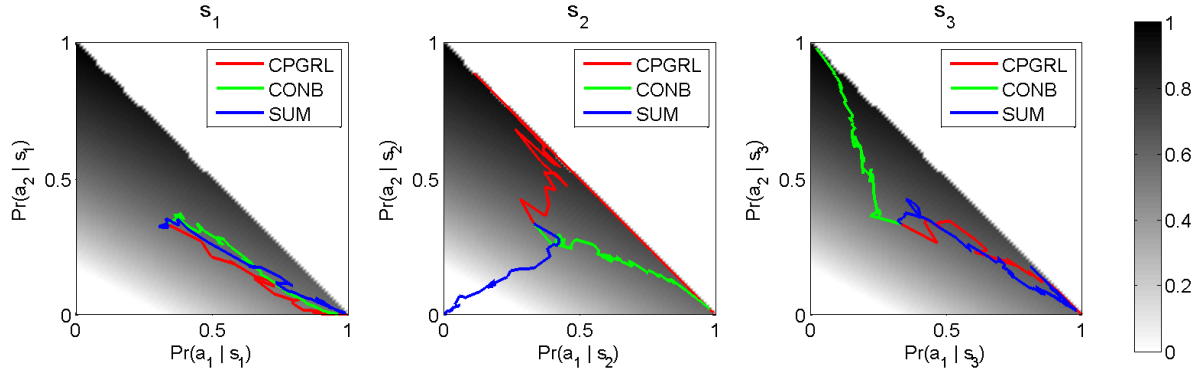


Fig. 4. Evolution of the policy parameters in the policy space.

The continuous state is $\mathbf{x} = [\phi_1, \phi_2]^\top$ while the continuous action consists of desired joint velocities, $\mathbf{u} = [d\phi_1, d\phi_2]^\top$. To represent the stochastic policy, we use a normalized Gaussian network,

$$\mu\theta(\mathbf{x}_t, \mathbf{u}_t) = \eta_1 \exp[-\eta_2 \|\mathbf{u}_t - \boldsymbol{\theta}^\top \mathbf{n}(\mathbf{x}_t)\|],$$

where η_1, η_2 and $\mathbf{n}(\mathbf{x})$ denote the constant values and the vector of the basis function. The number of basis functions are 40, determined by trial and error. In this experiment, (5) is used to compute the gradient.

This simulation does not consider dynamics of the arm. The initial joint angles are randomly initialized. The same meta-parameters such as learning rates are used in section III. When the arm touches one of the objects or $N_T = 1000$ time steps are expired, the episode terminates. One episode lasts for $N_K = 10000$ episodes and we perform 20 simulation runs.

B. Experimental results

Fig. 6 (a) shows the number of touches on the objects in each 100 episodes when the robotic arm is motivated only by the intrinsic reward. Since r^1 was a dense reward function, it

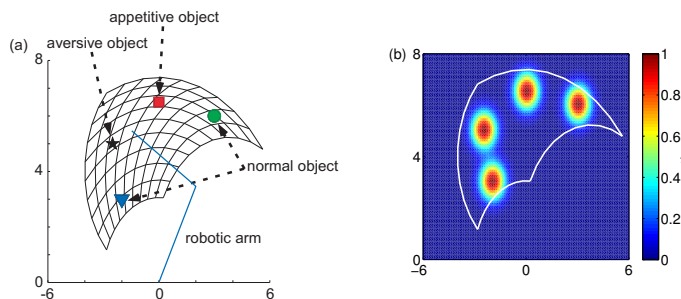


Fig. 5. Control task of a robotic arm. (a) Simulated environment where a mesh represents a reachable region of the end-effector of the arm. (b) Distribution of the intrinsic reward r^1 .

was not hard to obtain touching behaviors. Then we introduce two constraints by extrinsic rewards with thresholds $G^2 = G^3 = 0$. Fig. 6 (b) shows the experimental results. At the early stage of learning, the robotic arm touched the aversive star object. Then, it learned to avoid the aversive star object after about 1×10^3 episodes.

Finally, we strengthen the constraint by setting $G^3 = 0.5$ and observe the behaviors shown in Fig. 6 (c). Since the robotic arm has to touch the square object in order to obtain a positive r^2 , the number of touches on the square objects increases while those on other objects are gradually reduced to zero. It is revealed that G^2 is a sensitive threshold that affects the resultant behaviors.

V. DISCUSSION

In this paper, we have proposed the CPGRL that maximizes the long-term average reward under the inequality constraints. Experimental results encourage us to conduct the robotic experiments because one of our interests is to design the developmental learning methods for real hardware systems. Although we could not discuss the design principles of intrinsic and extrinsic rewards to establish a sustainable and scalable learning progress, this is very important. We think that the CPGRL gives the first step towards developmental learning. We develop the experimental setup that integrates the CPGRL and the technique of the embodied evolution in our multi-robot platform named *Cyber Rodents* [15]. In this case, the intrinsic reward is computed from sensor outputs while the extrinsic rewards are given according to the external events such as collisions with obstacles, capturing a battery pack, and so on. We show that good exploratory reward is acquired as the intrinsic reward through the interaction among three mobile robots [16]. We also plan to test other types of intrinsic rewards used in previous studies [3], [7].

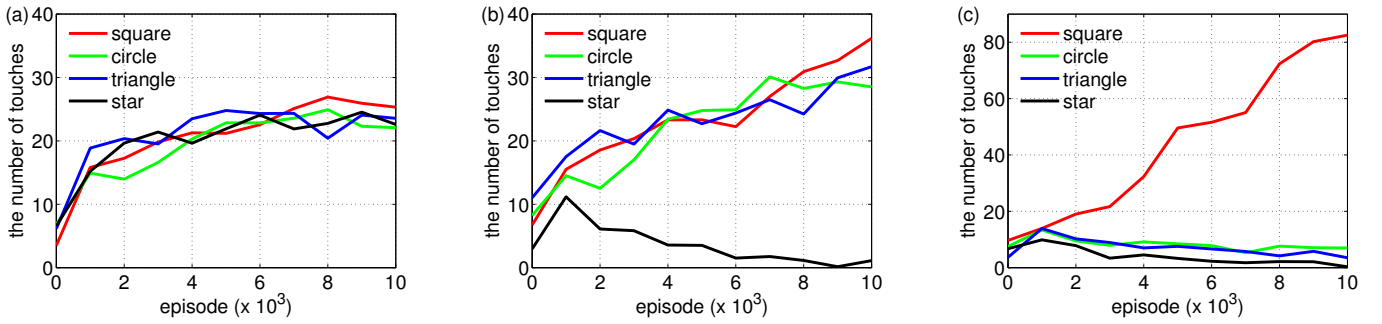


Fig. 6. Number of touches on the objects. (a) No constraints. (b) Normal constraints: $G^2 = G^3 = 0$. (c) Tight constraints: $G^2 = 0$ and $G^3 = 0.5$.

Finally, we describe three foreseeable extensions of this study. At first, we improve the efficiency of numerical computation. It is known that the learning speed of standard PGRL can be slow due to high variance in the estimate. Then, the Natural Policy Gradient (NPG) method [11] supported by the theory of information geometry is implemented to accelerate the speed of learning. Secondly, we develop a method to tune the thresholds used in the inequality constraints during learning processes. As shown in section IV-B, the learned behaviors were strongly affected by the setting of the thresholds. From a viewpoint of constrained optimization problems, G^i is just a meta-parameter given by the experimenters. However, the learning agent will show a variety of behaviors by changing these thresholds. We think that CPGRL has a potential to create new behaviors through the interaction between intrinsic and extrinsic rewards.

APPENDIX

Suppose that the policy parameter is modified without considering a restoration move d . When the k -th episode ends, the update rule is given by $\theta'_k = \theta_k + \alpha_1 s$, where s is the steepest ascent direction. By using the estimated gradients, the set of active constraints can be approximated by the following linear equation: $N_A^\top \theta + b \approx 0$ where b is an appropriate vector. Since the gradient projection method [12] assumes that θ lies in the subspace tangent to the active constraints, both θ'_k and θ_k should satisfy the above equations. Then, we obtain an equality constraints $N^\top s = 0$. Since the steepest ascent direction s satisfying the above constraints is required, we can pose the problem as

$$\max s^\top \Delta^1 \quad \text{s.t.} \quad N_A^\top s = 0 \quad \text{and} \quad s^\top s = 1. \quad (10)$$

In order to solve this problem with equality constraints, the Lagrangian function is defined by $\mathcal{L}(s, \lambda, \kappa) = s^\top \Delta^1 - s^\top N_A \lambda - \kappa (s^\top s - 1)$, where λ and κ are Lagrange multipliers. The condition for \mathcal{L} to be stationary is $\partial \mathcal{L} / \partial s = 0$. By multiplying N_A^\top and using the constraints, we obtain $N_A^\top \Delta^1 - N_A^\top N_A \lambda = 0$. Equation (9) is derived from this equation. Then, s is given by

$$s = \frac{1}{2\kappa} \left[I - N_A (N_A^\top N_A)^{-1} N_A^\top \right] \Delta^1,$$

and P should be defined in (7).

REFERENCES

- [1] R. S. Sutton and A. G. Barto. *Reinforcement Learning*. MIT Press/Bradford Books, 1998.
- [2] A. G. Barto, S. Singh, and N. Chentanez. Intrinsically Motivated Learning of Hierarchical Collections of Skills. In *Proc. of International Conference on Developmental Learning*, 2004.
- [3] S. Singh, A. G. Barto, and N. Chentanez. Intrinsically motivated reinforcement learning. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems 17*, pages 1281–1288. MIT Press, 2005.
- [4] A. Stout, G. D. Konidaris, and A. G. Barto. Intrinsically motivated reinforcement learning: A promising framework for developmental robot learning. In *Proc. of the AAAI Spring Symposium Workshop on Developmental Robotics*, 2005.
- [5] R. S. Sutton, D. Precup, and S. Singh. Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112:181–211, 1999.
- [6] L. A. Meeden, J. B. Marshall, and D. Blank. Self-Motivated, Task-Independent Reinforcement Learning for Robots. In *2004 AAAI Fall Symposium on Real-World Reinforcement Learning*, 2004.
- [7] P.-Y. Oudeyer and F. Kaplan. Intelligent adaptive curiosity: A source of self-development. In *Proc. of the 4th International Workshop on Epigenetic Robotics*, pages 127–130, 2004.
- [8] E. L. Deci and R. Flaste. *Why we do what we do: understanding self-motivation*. Penguin books, 1996.
- [9] J. Baxter and P. L. Bartlett. Infinite-horizon gradient-based policy search. *Journal of Artificial Intelligence Research*, 15(319–350), 2001.
- [10] V. R. Konda and J. N. Tsitsiklis. Actor-critic algorithms. *SIAM Journal on Control and Optimization*, 42(4):1143–1166, 2003.
- [11] T. Morimura, E. Uchibe, and K. Doya. Utilizing the natural gradient in temporal difference reinforcement learning with eligibility traces. In *Proc. of the 2nd International Symposium on Information Geometry and its Application*, pages 256–263, 2005.
- [12] S. A. Rosen. The gradient projection method for nonlinear programming — part I: linear constraints. *Journal of the Society for Industrial and Applied Mathematics*, 8(1):181–217, 1960.
- [13] E. Feinberg and A. Shwartz. Constrained dynamic programming with two discount factors: Applications and an algorithm. *IEEE Transactions on Automatic Control*, 44:628–630, 1999.
- [14] D. Dolgov and E. Durfee. Stationary deterministic policies for constrained MDPs with multiple rewards, costs, and discount factors. In *Proc. of the 19th International Joint Conference on Artificial Intelligence*, pages 1326–1331, 2005.
- [15] K. Doya and E. Uchibe. The Cyber Rodent Project: Exploration of adaptive mechanisms for self-preservation and self-reproduction. *Adaptive Behavior*, 13:149–160, 2005.
- [16] E. Uchibe and K. Doya. The cyber rodent project: Learning and evolution under the biological constraints. In *Proc. of the 14th International Conference on Neural Information Processing*, 2007. (submitted).