

Master's Thesis

Deep Reinforcement Learning for Optimal Self-Triggered Control

Guidance

Professor Yoshito OHTA
Assistant Professor Kenji KASHIMA

Ibuki TAKEUCHI

Department of Applied Mathematics and Physics

Graduate School of Informatics

Kyoto University



February 2021

Abstract

In this thesis, we construct the mathematical models of students who write the master's thesis and develop efficient algorithms for writing the master's thesis based on the models. We show that the proposed algorithms generate the thesis 65536 times more efficiently than writing by oneself.

Contents

1	Introduction	1
2	Preliminaries	1
2.1	Background of Deep Reinforcement Learning	1
2.2	Policy Iteration	1
2.3	Algorithms adapted to the settings of state and action space	2
2.4	Deterministic Policy Gradient Method	2
2.5	Exploration in DDPG	3
3	Problem Formulation	3
3.1	Self-Triggered Control	3
3.2	Optimal Self-Triggered Control	4
4	Reinforcement Learning for Optimal Self-Triggered Control	4
4.1	Model Settings	4
5	Consideration	5
5.1	Linear Case	5
5.1.1	Initial Policy	5
5.1.2	Learned Policy	6
5.1.3	Comparison with Naive Model Based Control	6
5.2	Non-Linear Case	6
5.2.1	Initial Policy	6
5.2.2	Learned Policy	6
5.3	Optimality	7
5.4	Approximation Accuracy of Critic	8
5.5	Ingenuity	8
6	Conclusion	9
	References	9
A	Appendix	10

1 Introduction

2 Preliminaries

2.1 Background of Deep Reinforcement Learning

Consider a malkov decision process M given with tuple $M = \{S, A, T, d_0, r, \gamma\}$. Here, S, A denotes state, action set, and $T(s'|s, a)$ express transition probability. Also, $d_0, r(s, a), \gamma \in [0, 1]$ are distribution of initial state, reward, discount factor respectively.

Now, the purpose of reinforcement learning is to find a policy such that

$$\pi^* = \operatorname{argmax}_{\pi} J(\pi) \quad (1)$$

where evaluation function $J(\pi)$ and (state) value function $V^\pi(s)$ is given as following:

$$V^\pi(s) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) |_{a_t=\pi(s_t)}, s_0 = s \quad (2)$$

$$J(\pi) = \mathbb{E}_{s_0 \sim d_0} [V^\pi(s_0)] \quad (3)$$

Here, we define Q function, which is useful tool for analyzing reinforcement learning. Q function is given as

$$\begin{aligned} Q^\pi(s, a) &= r(s, a) + \gamma \sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) |_{a_t=\pi(s_t)} \\ &= r(s, a) + \gamma V^\pi(s'). \end{aligned} \quad (4)$$

As showed in (4), Q function express the value when agent select action a freely and choose action according to the policy π from next step. Thus, the Q -function is also known as the action value function.

2.2 Policy Iteration

There is an algorithm for achieving (1), called the policy iteration method. It consists in repeating the following two steps.

1. Policy Evaluation: Find (or approximate) action value function $Q^\pi(s, a)$.
2. Policy Improvement: Update policy as $\pi(s) = \operatorname{argmax}_a Q^\pi(s, a)$.

It is known that the optimal policy π^* can be obtained by repeating the above two steps (Policy Improvement Theorem).

2.3 Algorithms adapted to the settings of state and action space

In the case that both the state space and the action space take discrete values, it is easy to obtain $\pi(s) = \underset{a}{\operatorname{argmax}} Q^\pi(s, a)$, which appeared in the previous section, by storing $Q^\pi(s, a)$ in a table.

Now, what about the case where the state space is continuous? Since the state s takes a continuous value, it cannot be stored in a table. Therefore, Minh et al.[2] took the approach of approximating $Q^\pi(s, a)$ by parametrizing it using a neural network. Since the action space is discrete, it is still possible to obtain $\underset{a}{\operatorname{argmax}} Q^\pi(s, a)$.

Finally, in the case where both state and action space is continuous, the problem is that it is very expensive to obtain $\underset{a}{\operatorname{argmax}} Q^\pi(s, a)$. Thus, up to this point, the policy π has been determined by the Q-function, but this approach cannot be taken when both spaces are continuous. Therefore, the policy function is often parameterized as π_θ and the parameter θ is updated by gradient ascent.

2.4 Deterministic Policy Gradient Method

Silver et al.[3] finds the gradient for the evaluation function $J(\pi_\theta)$, even if the policy $\pi(s)$ is defined as deterministic policy. This gradient is known as deterministic policy gradient(DPG), and it is calculate as following:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{s \sim \rho^{\pi_\theta}} [\nabla_\theta \pi_\theta(s) \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s)}] \quad (5)$$

where,

$$\rho^{\pi_\theta}(s) = \int_S \sum_{t=0}^{\infty} \gamma^t d_0(s_0) \Pr(s_0 \rightarrow s, t, \pi_\theta) ds_0 \quad (6)$$

is discounted distribution.

DDPG(Deep DPG) is a deep reinforcement learning algorithm which utilize this policy gradient. It adopts an Actor-Critic structure, and learns a critic network $Q(s, a|\omega)$ which approximates Q^{π_θ} , and an actor network $\pi(s|\theta) = \pi_\theta$ which represents a measure π , respectively. The update algorithm of actor and critic is described below.

DDPG uses mini-batch learning. First, update idea of critic is shown. The purpose of critic is to approximate Q^π . Because Q function is able to be decomposed like (4), $Q(s, a|\omega)$ should also be updated to satisfy this relation. For that, parameter ω is updated to the direction where it minimize Temporal Difference(TD) error:

$$\text{TD} = Q(s, a|\omega) - \{r(s, a) + \gamma Q(s, \pi(s)|\omega)\} \quad (7)$$

Since it is difficult to optimize for whole (s, a) at once, DDPG uses the mean squared error for the mini-batch E as the loss function and reduces it.

Now, the above method of updating critic is supervised learning itself. Therefore, the data in the mini-batch must be i.i.d.. If the mini-batch E uses the data of the last N steps experienced by the agent, they are no longer independent. Hence, the agent stores

the empirical data in a storage, called replay buffer, and randomly selects N data from it to make a mini-batch to increase the variance of the mini-batch.

Next update law of actor are shown. Because actor is the representation of policy function $\pi(s)$, policy gradient is used for its update. However, since correct Q -function as in equation (5) cannot be used in DDPG, approximated policy gradient

$$\mathbb{E}_{s \sim \rho^\pi} [\nabla_\theta \pi_\theta(s) \nabla_a Q(s, a | \omega)|_{a=\pi_\theta(s)}] \simeq \nabla_\theta J(\pi_\theta) \quad (8)$$

is used. Furthermore, an approximation to the expectation is made like following:

$$\mathbb{E}_{s \sim \rho^\pi} [\nabla_\theta \pi_\theta(s) \nabla_a Q(s, a | \omega)|_{a=\pi_\theta(s)}] \simeq \frac{1}{N} \sum_{s \in E} [\nabla_\theta \pi_\theta(s) \nabla_a Q(s, a | \omega)|_{a=\pi_\theta(s)}]. \quad (9)$$

Therefore, the accuracy of the approximation of the policy gradient may be greatly degraded by the accuracy of critic approximation and the distribution of mini-batch.

2.5 Exploration in DDPG

Reinforcement learning uses empirical data collected by interaction with the environment to improve a policy. An agent needs to take a variety of actions in each state in order to know better actions. This is called "exploration". When the action space is discrete, a exploration method called " ϵ -greedy" is widely known.

In the case of continuous action space, it is common to collect data with the action determined by adding noise e to the output of the policy function during training, such as

$$a = \pi_\theta(s) + e. \quad (10)$$

(The noise-added policy is called the behavior policy in order to distinguish it from the learning policy.) Intuitively, the larger this noise is, the wider the exploration is expected to be. By increasing the distribution of behavioral data around $a = \pi_\theta(s)$, we can expect to improve the accuracy of approximating the gradient of the Q -function with respect to a .

Then, is it useful to enlarge the exploration noise in all cases? In DDPG, we stored the past empirical data in the replay buffer, and calculated the policy gradient using the data set selected with equal probability from them. Since the expected value of the policy gradient is taken for $s \sim \rho^{\pi_\theta}$, the state distribution of the replay buffer should be as close as possible to ρ^{π_θ} . Adding a large exploration noise may violate this property in some cases. Therefore, it may be disadvantageous to increase the exploration noise. This is considered to be a kind of problem called the "exploration and exploitation dilemma".

3 Problem Formulation

3.1 Self-Triggered Control

We consider control system like Fig. 1. Here, the control target is a continuous-time

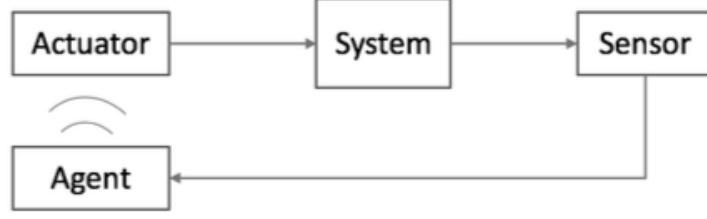


Figure 1: control system

system as in Equation (10)

$$\dot{s} = h(s, u) \quad (11)$$

Now, we consider a feedback control for system (10). In this paper, we call the observation of the state variable s and the sending of the input signal to the actuator as "interaction". In the self-triggered control, the agent does not make interaction continuously, but after a communication interval τ seconds determined by the agent itself. In order to express it mathematically, we assume that the agent's control law $\pi(s)$ is a vector-valued function consisting of two elements, where the first element represents the input u sent to the actuator and the second element represents the interval $\tau(sec)$. The input u sent in the previous interaction is added until the time of the next interaction (ZOH control).

3.2 Optimal Self-Triggered Control

The optimal self-triggered control law π^* is defined as follows

$$\pi^* = \underset{\pi}{\operatorname{argmax}} J(\pi) \quad (12)$$

$$J(\pi) = \mathbb{E}_{s_0 \in d_0} [V^\pi(s_0)] \quad (13)$$

$$V^\pi(s_0) = \sum_{i=0}^{\infty} \gamma^i r_i^\pi \quad (14)$$

$$r_i^\pi = - \int_{T_i}^{T_{i+1}} s(t)^\top Q s(t) dt + \tau_i a_i^\top R a_i + \lambda \tau_i, \quad T_i = \sum_{l=0}^i \tau_l \quad (15)$$

where i denote the number of interactions, and a_i and τ_i be the outputs of the measures π for the i th interaction, respectively.

4 Reinforcement Learning for Optimal Self-Triggered Control

In reinforcement learning, the agent is expected to perform interactions step by step. The self-triggered control problem considered in this paper, where the control target is

given as a continuous-time system, satisfies the assumption of reinforcement learning by considering the communication with the actuator as one step.

4.1 Model Settings

In the self-triggered control, the agent needs to decide the input signal u and the interval τ at each step. In this paper, the control law is given as a state feedback. Therefore, the policy function is given as follows:

$$\pi(s) = [u(s) \quad \tau(s)]^\top \quad (16)$$

We use DDPG as reinforcement learning algorithm. As described in section 2, actor and critic is expressed as neural networks respectively. Experiments have shown that learning diverges when using a general network, so here we use the special network. The architecture of 2 networks is in Fig 2.

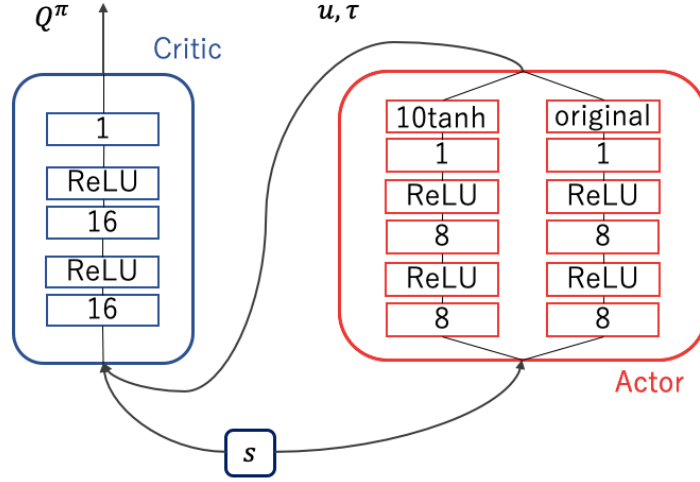


Figure 2: Agent Model

The activation function "original" shown in Fig 2 is defined as $0.99 \times \text{sigmoid} + 0.01$ to meet upper and lower limits of interval.

5 Consideration

In this section, we study the effectiveness of the reinforcement learning approach to the optimal self-triggered control problem. We conduct numerical experiments and review the results for the cases of linear and nonlinear control systems, respectively. In both cases, the communication interval is allowed to be $0.01(s) \sim 1.0(s)$.

5.1 Linear Case

First, we adopt reinforcement learning to self-triggered control for linear system. The control object is

$$\dot{s} = As + Bu = \begin{bmatrix} -1 & 4 \\ 2 & -3 \end{bmatrix} s + \begin{bmatrix} 2 \\ 4 \end{bmatrix} u. \quad (17)$$

Here, the input signal u is limited to $-10 \sim 10$.

5.1.1 Initial Policy

For learning efficiency, we use π_{init} such that

$$\pi_{\text{init}} = [-Ks \quad 0.01] \quad (18)$$

where K is a feedback gain calculated by Linear Quadratic Regulator. This policy stabilize the system (really??).

5.1.2 Learned Policy

5.1.3 Comparison with Naïve Model Based Control

5.2 Non-Linear Case

In this subsection, we investigate whether the self-triggered control law can be learned by reinforcement learning even when the control target is extended to non-linear systems, especially control affine systems. We consider an inverted pendulum, whose state-space representation is

$$\frac{d}{dt} \begin{pmatrix} \theta \\ \dot{\theta} \end{pmatrix} = \begin{pmatrix} \dot{\theta} \\ \frac{3g}{2l} \sin \theta + \frac{3}{ml^2} a \end{pmatrix}. \quad (19)$$

Therefore, for an inverted pendulum, the state variable s is considered to be $(\theta \quad \dot{\theta})^\top$.

As in the linear case, the input signal u is limited to $-10 \sim 10$.

5.2.1 Initial Policy

For learning efficiency, we use π_{init} such that

$$\pi_{\text{init}} = [-Ks \quad 0.2] \quad (20)$$

where K is a feedback gain calculated by Linear Quadratic Regulator. This policy stabilize the system.

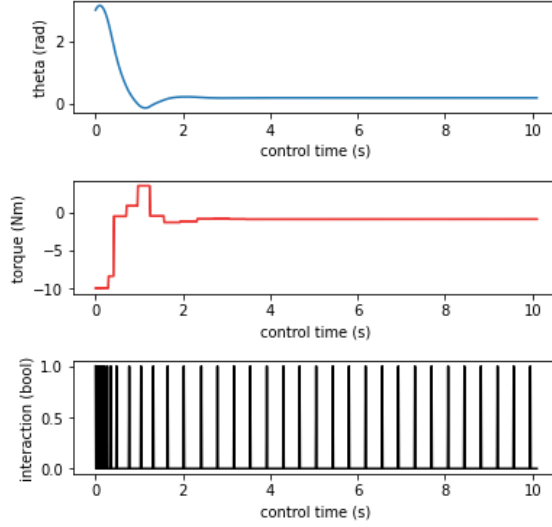


Figure 3: A control path with learned policy π_{RL}

5.2.2 Learned Policy

Fig 3 shows a control path with learned policy, stating from $s_0 = [3., 3.]$.

Figure 3 shows, from top to bottom, the angle of the pendulum θ , the control input u , and the boolean representing interaction. From this figure, we can see that the communication interval is determined flexibly at each interaction.

Now, we confirm whether the learned policy π_{RL} enlarge the evaluation function $J(\pi)$ compared from the initial policy π_{init} . Since $J(\pi) = \mathbb{E}_{s_0}[V^\pi(s_0)]$, we generate 500 initial states s_0 according to the initial state distribution d_0 , and average the results of $V^\pi(s_0)$. We conduct this approximation of expectation 1000 times. From the result

$$J(\pi_{\text{init}}) = 4.23 \pm 0.083, \quad J(\pi_{\text{RL}}) = 37.12 \pm 0.071, \quad (21)$$

we can confirm $J(\pi_{\text{init}}) < J(\pi_{\text{RL}})$. It can be concluded that the policy has been improved.

5.3 Optimality

Is the policy π_{RL} described in the previous section an optimal? If it is an optimal solution, the policy gradient must be $\mathbf{0}$. Therefore from equation (5), a condition

$$\nabla_a Q^{\pi_{\text{RL}}}(s, a)|_{a=\pi_{\text{RL}}(s)} = \mathbf{0} \quad (22)$$

should be met.

In order to check whether the policy π_{RL} satisfies the condition (??), we want to know $Q^{\pi_{\text{RL}}}(s, a)$. Then, we obtain the approximation of value function $V^{\pi_{\text{RL}}}(s) =$

$\sum_{i=0}^N \gamma^i r(s_i, \pi_{\text{RL}}(s_i))$ with N step simulation:

$$V_{\text{ap}}^{\pi_{\text{RL}}}(s) = \sum_{i=0}^N \gamma^i r(s_i, \pi_{\text{RL}}(s_i)). \quad (23)$$

Next, the approximation of $Q^{\pi_{\text{RL}}}$ is calculated with $V_{\text{ap}}^{\pi_{\text{RL}}}$ as following:

$$Q_{\text{apd}}^{\pi_{\text{RL}}}(s, a) = r(s, a) + \gamma V_{\text{ap}}^{\pi_{\text{RL}}}(s') \quad (24)$$

Since $\gamma < 1$, $Q_{\text{apd}}^{\pi_{\text{RL}}}(s, a)$ is a good approximation if the number of simulation steps N is sufficiently large. However, to compute the policy gradient, we need to compute $\nabla_a Q^{\pi_{\text{RL}}}(s, a)$, while $Q_{\text{apd}}^{\pi_{\text{RL}}}$ can only be computed the value at (s, a) . Then, we approximate the function $Q_{\text{apd}}^{\pi_{\text{RL}}}(s, a)$ by supervised learning with $Q_{\text{apd}}^{\pi_{\{extrm{RL}\}}}(s, a)$ as the teacher data. This function is defined as $Q_{\text{ap}}^{\pi_{\text{RL}}}$.

Now, when we calculate the policy gradient using $Q_{\text{ap}}^{\pi_{\text{RL}}}$, we obtain

$$\mathbb{E}_{s \sim \rho^{\pi_{\text{RL}}}} [\nabla_{\theta} \pi_{\text{RL}}(s) \nabla_a Q^{\pi_{\text{RL}}}(s, a)|_{a=\pi_{\text{RL}}(s)}] \neq \mathbf{0}, \quad (25)$$

which means that the policy π_{RL} is not an optimal measure. We will discuss the reason for this.

5.4 Approximation Accuracy of Critic

Since policy π_{RL} makes

$$\mathbb{E}_{s \sim \rho^{\pi_{\text{RL}}}} [\nabla_{\theta} \pi_{\text{RL}}(s) \nabla_a Q(s, a|\omega)|_{a=\pi_{\text{RL}}(s)}] \simeq \mathbf{0} \quad (26)$$

for critic network $Q(s, a|\omega)$, it is a optimal solution to make the approximated gradient be 0. Therefore, the reason why we could not learn the optimal strategy is considered to be the low approximation accuracy of critic. In this experiment, the exploration noise is given as \sim , and its scale is small. This leads to a bias in the distribution of the action a in the replay buffer. For example, the distribution of actions experienced in the neighborhood of $s = [0 \ 0]$ is shown in Figure 4.

Since the approximation accuracy of

$$\nabla_a Q(s, a|\omega)|_{a=\pi_{\theta}(s)} \quad (27)$$

is crucial for the calculation of the policy gradient, the problem is the lack of data around $a = \pi_{\theta}(s)$ which is represented by the orange dot in Figure 4.

5.5 Ingenuity

We see that small scale exploration noise leads low accuracy of critic's approximation. However, as we confirmed in section 2.5, large scale exploration noise cause problem of experience distribution. Here, by using the information of the control target, we consider to devise a exploration method.

The ideal distribution of replay buffer meets following 2 conditions:

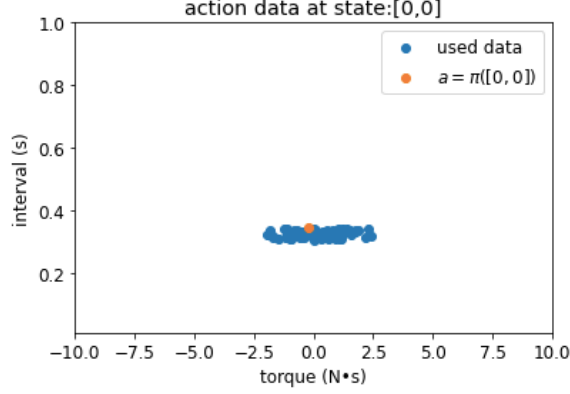


Figure 4: Action distribution in replay buffer around $s = [0 \ 0]$.

1. The distribution of state s is similar to discounted distribution $\rho^\pi(s)$
2. The variance of action a is large for each state s .

Therefore, we propose a method to control the noise scale inversely proportional to the gradient of the state s' of the next step with respect to the change of input a .

$$e \sim \frac{1}{c\|g\| + 1} \mathcal{N}(0, 1) \quad (28)$$

where $g = \frac{ds'}{da}$, c is a hyper parameter.

As a result of reinforcement learning using this search method, the obtained policy π_{MBRL} achieve

$$J(\pi_{\text{MBRL}}) = 68.10 \pm 0.018. \quad (29)$$

6 Conclusion

本稿では, 最適セルフトリガー制御に対する強化学習

Acknowledgments

The author would like to express his sincere gratitude to Professor Yoshito Ohta and Assistant Professor Kenji Kashima for their helpful advices.

References

- [1] C. J. Watkins, and P. Dayan. Q-learning. *Machine Learning*, vol. 8, no. 3-4, pp. 279-292, 1992.

- [2] V. Minh, K. Kavukcoglu, D. Silver. et al.. “Human-level control through deep reinforcement learning.” *Nature* 518, pp.529-533, 2015.
- [3] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, et al.. “Deterministic Policy Gradient Algorithms.” *ICML Beijing, China.*, 2014, Beijing.
- [4] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N.Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *International Conference on Learning Representations*, 2015.
- [5] T. Degris, M. White and R. Sutton. “Off-Policy Actor-Critic.” *ICML Edinburgh, United Kingdom*, 2012.
- [6] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. ”Policy gradient methods for reinforcement learning with function approximation.” *In Advances in Neural Information Processing Systems*, 2000.
- [7] D. P. Kingma and J. Ba. “Adam: A Method for Stochastic Optimization.” *arXiv preprint arXiv: 1412.6980*, 2014.

A Appendix

This is an appendix. This is a citation [?].

Table 1: This is a table.

	A	B
C	70	80
D	100	0

Master's Thesis

Deep Reinforcement Learning for Optimal Self-Triggered Control

Guidance

Professor Yoshito OHTA
Assistant Professor Kenji KASHIMA

Ibuki TAKEUCHI

Department of Applied Mathematics and Physics

Graduate School of Informatics

Kyoto University



February 2021

Deep Reinforcement Learning for Optimal Self-Triggered Control

Ibuki TAKEUCHI

February 2021

Deep Reinforcement Learning for Optimal Self-Triggered Control

Ibuki TAKEUCHI

Abstract

In this thesis, we construct the mathematical models of students who write the master's thesis and develop efficient algorithms for writing the master's thesis based on the models. We show that the proposed algorithms generate the thesis 65536 times more efficiently than writing by oneself.