

Master's Thesis

Deep Reinforcement Learning for 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	2
2.1	Background of Deep Reinforcement Learning	2
2.2	Policy Iteration	3
2.3	Algorithms adapted to the settings of state and action space	3
2.4	Deterministic Policy Gradient Method	3
3	Problem Formulation	5
3.1	Self-Triggered Control	5
3.2	Previous Research for Self-Triggered Control	6
3.3	Optimal Self-Triggered Control	6
4	Reinforcement Learning for Self-Triggered Control	7
4.1	Deterministic Policy Gradient for Self-Triggered Control	7
4.2	Value function for Self-Triggered Control	9
4.3	Ideal Algorithm	10
4.4	Practical Algorithm	11
4.5	Other approach	11
5	Consideration	12
5.1	Linear System	12
5.1.1	initial policy	12
5.1.2	proposed method 1	13
5.1.3	proposed method 2	13
5.1.4	Comparison with Naïve Model Based Control	14
5.2	Non-Linear Case	14
5.2.1	Initial Policy	14
5.2.2	proposed method 1	14
5.2.3	proposed method 2	14
6	Conclusion	15
	References	15
A	Appendix	16
A.1	Model Settings	16

1 Introduction

In many control applications, controllers are nowadays implemented using communication networks in which the control task has to share the communication resources with other tasks. Despite the fact that resources can be scarce, controllers are typically still implemented in a time-triggered fashion, in which control tasks are executed periodically. This design choice often leads to over-utilization of the available communication resources, and/or causes a limited lifetime of battery-powered devices, as it might not be necessary to execute the control task every period to guarantee the desired closed-loop performance.

Also in the area of "sparse control" (Gallieri & Maciejowski, 2012), in which it is desirable to limit the changes in certain actuator signals while still realizing specific control objectives, periodic execution of control tasks may not be optimal either. In both networked control systems with scarce communication resources and sparse control applications arises the fundamental problem of determining optimal sampling and communication strategies, where optimality needs to reflect both implementation cost (related to the number of communications and/or actuator changes) as well as control performance. It is expected that the solution to this problem results in control strategies that abandon the periodic time-triggered control paradigm.

Two approaches that abandon the periodic communication pattern are event-triggered control (ETC), see, e.g., Arzen (1999), Astrom and Bernhardsson (1999) and Donkers and Heemels (2012), Heemels, Sandee, and van den Bosch (2008), Heemels et al. (1999), Henningsson, Johansson, and Cervin (2008), Lunze and Lehmann (2010), Tabuada (2007) and Wang and Lemmon (2009), and self-triggered control (STC), see, e.g., Almeida, Silvestre, and Pascoal (2010, 2011), Anta and Tabuada (2010), Donkers, Tabuada, and Heemels (2012), Mazo, Anta, and Tabuada (2010), Velasco, Fuertes, and Marti (2003) and Wang and Lemmon (2009). Although ETC is effective in the reduction of communication or actuator movements, it was originally proposed for different reasons, including the reduction of the use of computational resources and dealing with the event-based nature of the plants to be controlled. In ETC and STC, the control law consists of two elements being a feedback controller that computes the control input, and a triggering mechanism that determines when the control input has to be updated. The difference between ETC and STC is that in the former the triggering consists of verifying a specific condition continuously and when it becomes true, the control task is triggered, while in the latter at an update time the next update time is pre-computed.

At present ETC and STC form popular research areas. However, two important issues have only received marginal attention: (i) the co-design of both the feedback law and the triggering mechanism, and (ii) the provision of performance guarantees (by design). To elaborate on (i), note that current design methods for ETC and STC are mostly emulation-based approaches, by which we mean that the feedback controller is designed without considering the scarcity in the system's resources. The triggering mechanism is only designed in a subsequent phase, where the controller has already been fixed. Since the feedback controller is designed before the triggering mechanism, it is difficult,

if not impossible, to obtain an optimal design of the combined feedback controller and triggering mechanism in the sense that the minimum number of control executions is achieved while guaranteeing closed-loop stability and a desired level of performance.

By the way, artificial intelligence is nowadays used in various situations, notably in automatic driving technology, and the development of research on the subject of artificial intelligence is remarkable. One of the concepts to realize artificial intelligence is reinforcement learning. Reinforcement learning is an algorithm that learns behaviors that maximize the long-term benefits by repeated trial and error. In addition, although not mathematically proven, reinforcement learning has been used to obtain meaningful results for nonlinear systems. In this paper, we investigate the usefulness of reinforcement learning as a method to realize the self-triggered control law.

The two main contributions of this research are

- To formulate the optimal self-triggered control problem for long-time costs explicitly considering communication to costs, and to consider the policy gradient for it.
- To obtain a state feedback control law that outperforms the control performance of a naively designed simulation-based self-triggered control law.
- To confirm the usefulness of reinforcement learning for self-triggered control not only for linear systems but also for non-linear systems.

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 $\pi(s) = \operatorname{argmax}_a Q^\pi(s, a)$ 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 $\operatorname{argmax}_a Q^\pi(s, a)$.

Finally, in the case where both state and action space is continuous, the problem is that it is very expensive to obtain $\operatorname{argmax}_a Q^\pi(s, a)$. 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 method.

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 theorem.

Theorem 1 (Deterministic Policy Gradient Theorem). *The gradient for evaluation function $\nabla_\theta J(\pi_\theta)$ is exist and calculated as,*

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

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.

Proof. First, we consider the gradient for $V^{\pi_\theta}(s)$.

$$\begin{aligned}
& \nabla_\theta V^{\pi_\theta}(s) \\
&= \nabla_\theta Q^{\pi_\theta}(s, \pi_\theta(s)) \\
&= \nabla_\theta [r(s, \pi_\theta(s)) + \gamma \int_S Pr(s \rightarrow s', 1, \pi_\theta) V^{\pi_\theta}(s') ds'] \\
&= \nabla_\theta \pi_\theta(s) \nabla_a r(s, a)|_{a=\pi_\theta(s)} \\
&\quad + \gamma \int_S (\nabla_\theta \pi_\theta(s) \nabla_a Pr(s \rightarrow s', 1, a)|_{a=\pi_\theta(s)} V^{\pi_\theta}(s') \\
&\quad + Pr(s \rightarrow s', 1, \pi_\theta) \nabla_\theta V^{\pi_\theta}(s')) ds' \\
&= \nabla_\theta \pi_\theta(s) \nabla_a [r(s, a) + \gamma \int_S Pr(s \rightarrow s', 1, \pi_\theta) V^{\pi_\theta}(s')]_{a=\pi_\theta(s)} ds' \\
&\quad + \gamma \int_S Pr(s \rightarrow s', 1, \pi_\theta) \nabla_\theta V^{\pi_\theta}(s') ds' \\
&= \nabla_\theta \pi_\theta(s) \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s)} + \gamma \int_S Pr(s \rightarrow s', 1, \pi_\theta) \nabla_\theta V^{\pi_\theta}(s') ds' \tag{7}
\end{aligned}$$

By using this relation recursively, we have,

$$\begin{aligned}
\nabla_\theta V^{\pi_\theta}(s) &= \sum_{i=0}^{\infty} \int_S \cdots \int_S Pr(s \rightarrow s', 1, \pi_\theta) Pr(s' \rightarrow s'', 1, \pi_\theta) \cdots \\
&\quad \gamma^i \nabla_\theta \pi_\theta(s' \cdots i) \nabla_a Q^{\pi_\theta}(s' \cdots i, a)|_{a=\pi_\theta(s' \cdots i)} ds' \cdots ds' \\
&= \sum_{i=0}^{\infty} \gamma^i \int_S Pr(s \rightarrow s', i, \pi_\theta) \nabla_\theta \pi_\theta(s) \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s')} ds'. \tag{8}
\end{aligned}$$

Since $J(\pi) = \mathbb{E}_{s \sim d_0}[V^\pi(s)]$,

$$\begin{aligned}
\nabla_\theta J(\pi_\theta) &= \nabla_\theta \int_S d_0(s) V^{\pi_\theta}(s) ds \\
&= \int_S d_0(s) \nabla_\theta V^{\pi_\theta}(s) ds \\
&= \int_S \rho^{\pi_\theta} \nabla_\theta \pi_\theta(s) \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s)} ds \tag{9}
\end{aligned}$$

□

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:

$$\text{TD} = Q(s, a|\omega) - \{r(s, a) + \gamma Q(s, \pi(s)|\omega)\} \quad (10)$$

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.

$$\text{Loss} = \frac{1}{N} \sum_{s \in E} \text{TD}^2 \quad (11)$$

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 it.

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

$$\mathbb{E}_{s \sim \rho^\pi} [\nabla_\theta \pi_\theta(s) \nabla_a Q(s, a|\omega)|_{a=\pi_\theta(s)}] \simeq \nabla_\theta J(\pi_\theta) \quad (12)$$

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

$$\mathbb{E}_{s \sim \rho^\pi} [\nabla_\theta \pi_\theta(s) \nabla_a Q(s, a|\omega)|_{a=\pi_\theta(s)}] \simeq \frac{1}{N} \sum_{s \in E} [\nabla_\theta \pi_\theta(s) \nabla_a Q(s, a|\omega)|_{a=\pi_\theta(s)}]. \quad (13)$$

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.

3 Problem Formulation

3.1 Self-Triggered Control

We consider control system like Fig 1.

Here, the control target is a continuous-time system as Equation (14)

$$\dot{s} = h(s, u) + \dot{w} \quad (14)$$

Now, we consider a feedback control for system (14). 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 π 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).

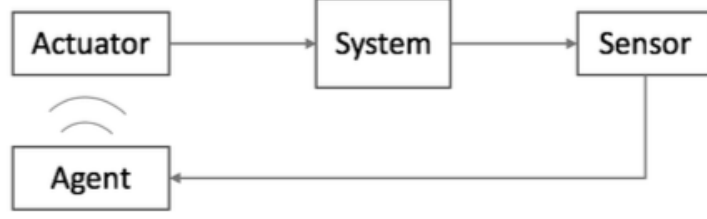


Figure 1: control system

3.2 Previous Research for Self-Triggered Control

Here, we summarize previous studies on self-triggered control. (write)

In our research, we formulate an optimal self-triggered control problem in which communication is explicitly included in the cost. This makes it possible to consider the problem of finding a control policy with a long-term cost, instead of a one-step optimization.

3.3 Optimal Self-Triggered Control

In self-triggered control, the agent needs to decide the input signal u and the interval τ at each step. Thus, the action a in reinforcement learning corresponds to $[u \ \tau]^\top$. (In this paper, we equate a with the tuple (u, τ) .)

Now, in this research, the control law is given as a state feedback. Therefore, the policy function is given as follows:

$$\pi(s) = [u(s) \ \tau(s)]^\top \quad (15)$$

In order to converge to the origin state as quickly as possible with the minimum input energy while reducing the frequency of communication, the agent aims to find a policy π^* that minimizes the following expected discounted cost

$$J(\pi) = \mathbb{E}_{s \sim d_0}[V^\pi(s)] \quad (16)$$

where

$$V^\pi(s) = \int_0^\infty e^{-\alpha t} \mathbb{E}_w[s(t)^\top Q s(t) + u(t)^\top R u(t) + \beta C(t) | s(0) = s] dt. \quad (17)$$

If we separate the definite integral for each interval, we have

$$\begin{aligned}
V^\pi(s) &= \sum_{i=0}^{\infty} \int_{t_i}^{t_{i+1}} e^{-\alpha t} \mathbb{E}_w[s(t)^\top Qs(t) + u_i^\top Ru_i + \beta C(t) | s(0) = s] dt \\
&= \sum_{i=0}^{\infty} e^{-\alpha t_i} \left(\int_0^{\tau_i} \mathbb{E}_w[s(t)^\top Qs(t) + u_i^\top Ru_i | s(0) = s_i] dt + \beta \right) \\
&= \sum_{i=0}^{\infty} e^{-\alpha t_i} r(s_i, \pi(s_i)).
\end{aligned} \tag{18}$$

Here, t_i is the time of the i th communication and s_i is the state at that time. Also, let $[u_i, \tau_i] = \pi(s_i)$. Therefore, $V^\pi(s)$ satisfies the following Bellman equation.

$$V^\pi(s) = r(s, \pi_\theta(s)) + e^{-\alpha \tau} \mathbb{E}_{s'}[V^\pi(s'(s, \pi_\theta(s)))] \tag{19}$$

In addition, the action value function Q^π is the discounted accumulation cost when agent freely chooses an action in the first step and follows the policy π from the next step, which satisfies the following Bellman equation.

$$Q^\pi(s, u, \tau) = r(s, u, \tau) + e^{-\alpha \tau} \mathbb{E}_{s'}[Q^\pi(s'(s, u, \tau), \pi(s'(s, u, \tau)))] \tag{20}$$

4 Reinforcement Learning for Self-Triggered Control

In this section, we consider the application of reinforcement learning to find the optimal self-trigger policy π^* . Simply thinking, we can consider the reinforcement learning problem by taking the communication as one step. Furthermore, DDPG may also be applied by approximating the Q -function, which satisfies the equation (20) using a critic network, to obtain the policy gradient. In this section, we discuss the validity of this approach.

4.1 Deterministic Policy Gradient for Self-Triggered Control

Since the discount factor in Equation (18) depends on τ at each step, it differs from the general reinforcement learning problem. In this subsection, we discuss how the DDPG is affected by this difference.

Actually, due to the property of Q -functions such as (20), DPG cannot be computed

as in (5). Since

$$\begin{aligned}
\nabla_\theta V^{\pi_\theta}(s) &= \nabla_\theta Q^{\pi_\theta}(s, \pi_\theta(s)) \\
&= \nabla_\theta[r(s, \pi_\theta(s)) + e^{-\alpha\tau_\theta(s)} \mathbb{E}_{s'}[V^{\pi_\theta}(s')]] \\
&= \nabla_\theta \pi_\theta(s) \nabla_a r(s, a)|_{a=\pi_\theta(s)} \\
&\quad + e^{-\alpha\tau_\theta(s)} \int_S \{ \nabla_\theta \pi_\theta(s) \nabla_a Pr(s \rightarrow s', 1, a)|_{a=\pi(s)} V^{\pi_\theta}(s') \\
&\quad \quad + Pr(s \rightarrow s', 1, \pi_\theta) \nabla_\theta V^{\pi_\theta}(s') \} ds' \\
&\quad + \int_S \nabla_\theta e^{-\alpha\tau_\theta(s)} Pr(s \rightarrow s', 1, \pi_\theta) V^{\pi_\theta}(s') ds', \tag{21}
\end{aligned}$$

we have

$$\begin{aligned}
\nabla_\theta V^{\pi_\theta}(s) &= \sum_{i=0}^{\infty} \int_S \cdots \int_S Pr(s_0 \rightarrow s_1, 1, \pi_\theta) \cdots Pr(s_{i-1} \rightarrow s_i, 1, \pi_\theta) \\
&\quad e^{-\alpha t_i} \nabla_\theta \pi_\theta(s_i) \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s_i)} ds_i ds_{i-1} \cdots ds_1 \\
&\quad + \sum_{i=1}^{\infty} \int_S \cdots \int_S Pr(s_0 \rightarrow s_1, 1, \pi_\theta) \cdots Pr(s_{i-1} \rightarrow s_i, 1, \pi_\theta) \\
&\quad e^{-\alpha t_{i-1}} \nabla_\theta e^{-\alpha\tau_\theta(s_{i-1})} V^{\pi_\theta}(s_i) ds_i ds_{i-1} \cdots ds_1 \\
&= \sum_{i=0}^{\infty} \int_S \cdots \int_S Pr(s_0 \rightarrow s_1, 1, \pi_\theta) \cdots Pr(s_{i-1} \rightarrow s_i, 1, \pi_\theta) \\
&\quad e^{-\alpha t_i} \nabla_\theta \pi_\theta(s_i) \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s_i)} ds_i ds_{i-1} \cdots ds_1 \\
&\quad + \sum_{i=0}^{\infty} \int_S \cdots \int_S Pr(s_0 \rightarrow s_1, 1, \pi_\theta) \cdots Pr(s_i \rightarrow s_{i+1}, 1, \pi_\theta) \\
&\quad e^{-\alpha t_i} \nabla_\theta e^{-\alpha\tau_\theta(s_i)} V^{\pi_\theta}(s_{i+1}) ds_{i+1} ds_i \cdots ds_1, \tag{22}
\end{aligned}$$

where

$$t_i = \begin{cases} 0 & (i = 0) \\ \sum_{k=0}^{i-1} \tau_\theta(s_{k+1}) & (otherwise) \end{cases}. \tag{23}$$

Thus, we can conclude that DPG for self-triggered control can be written as

$$\begin{aligned}
\nabla_\theta J(\pi_\theta) &= \mathbb{E}_{s_0 \sim d_0} [\nabla_\theta V^{\pi_\theta}(s_0)] \\
&= \sum_{i=0}^{\infty} \int_S \cdots \int_S d_0(s_0) Pr(s_0 \rightarrow s_1, 1, \pi_\theta) \cdots Pr(s_i \rightarrow s_{i+1}, 1, \pi_\theta) \\
&\quad e^{-\alpha t_i} \{ \nabla_\theta \pi_\theta(s_i) \nabla_a Q^{\pi_\theta}(s, a)|_{a=\pi_\theta(s_i)} + \nabla_\theta e^{-\alpha\tau_\theta(s_i)} V^{\pi_\theta}(s_{i+1}) \} ds_{i+1} ds_i \cdots ds_0. \tag{24}
\end{aligned}$$

4.2 Value function for Self-Triggered Control

In DPG for reinforcement learning problems with a fixed discount factor, it is sufficient that the gradient of the critic $Q(s, a|\omega)$ with respect to a can correctly approximate that of the Q -function. However, in the case of reinforcement learning for self-triggered control considered in this paper, the value of $Q(s, \pi(s)|\omega)$ itself must also be correctly approximated. Therefore, we need to pay attention to the TD learning of critic.

In this section, we discuss whether the critic learned using TD learning can approximate the value of $Q^\pi(s, a)$. First of all, since the Bellman equation for the Q -function in self-triggered control should satisfy (20), the critic should set the TD error

$$TD = Q(s, u, \tau|\omega) - \{r(s, u, \tau) + e^{-\alpha\tau} \mathbb{E}_{s'}[Q(s'(s, u, \tau), \pi(s'(s, u, \tau))|\omega)]\} \quad (25)$$

to be zero for all (s, u, τ) . In this section, we discuss the algorithm for learning such a critic.

Here, we create a mini-batch E from the dataset $D = (s, u, \tau, r, s')$ and learn critic by minimizing the MSE of the TD error for the mini-batch E . Therefore, if there is a bias of the distribution of (s, u, τ) in the empirical data set D , the accuracy of function approximation outside the distribution will obviously be low. From the equation (24), the approximation of $Q(s, \pi(s)|\omega)$ and $\nabla_a Q(s, a|\omega)|_{a=\pi(s)}$ is necessary for the calculation of the directional gradient, so we create a dataset D which contains s in the whole region of S and $[u, \tau] = \pi(s) + e$ (where e is a stochastic noise)

Algorithm 1 shows the learning algorithm for critics described in this section. In

Algorithm 1 TD learning for critic network

Sample M states s with equal probability from state space S .

for $r = 0$ to R **do**

For all sampled s , choose $[u, \tau] = \pi(s) + e$.

Execute action u for τ second to the environment.

Receive r and observe next state s'

Store (s, u, τ, r, s') to data set D

end for

for epoch = 0 to N **do**

Select m data pairs (s, u, τ, r, s') from D and make a mini-batch E .

Calculate gradient $g = \frac{\partial}{\partial \omega} \frac{1}{m} \sum (Q(s, u, \tau|\omega) - \{r + e^{-\alpha\tau} Q(s', \pi(s'))|\omega\})^2$

Update ω with gradient g

end for

the last part of this section, we compare $Q^\pi(s, \pi(s))$ for a self-triggered control law π with the critic $Q(s, \pi(s)|\omega)$ which approximates $Q^\pi(s, \pi(s))$ using the algorithm 1. Both $Q^\pi(s, \pi(s))$ and $Q(s, \pi(s)|\omega)$ are functions of state s . The state s is assumed to be two-dimensional, and the comparison between them is shown in 2.

In Figure 2, the blue points indicate the true Q^π and the orange points indicate the critics. The true Q^π is obtained by simulation. From Figure 2, we can see that the critic learned by Algorithm 1 is a good approximation of Q^π .

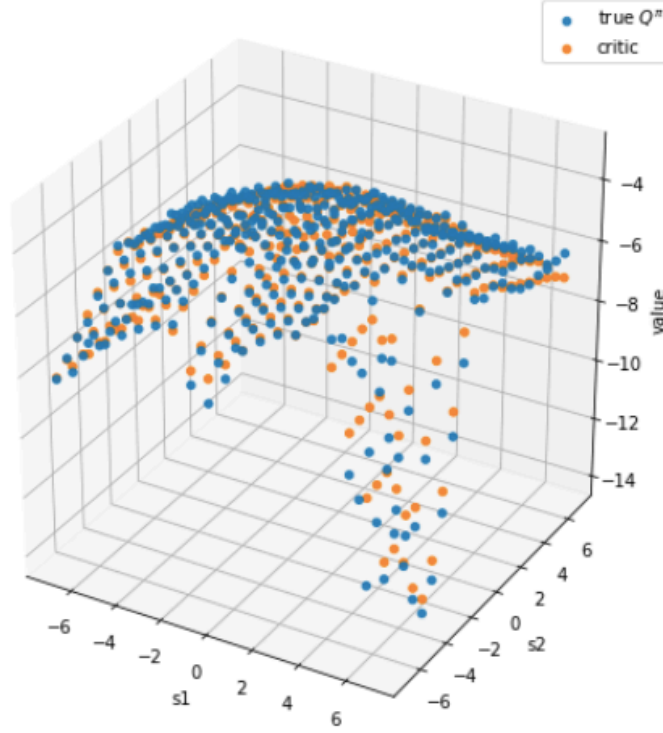


Figure 2: Approximation of $Q^\pi(s, \pi(s))$

4.3 Ideal Algorithm

We start by considering an ideal algorithm for computing the exact policy gradient (24) without considering the computational complexity. First, for an initial state s_0 , we perform P episodes of control for T seconds. In each step, we store $(s_i, u_i, \tau_i, r_i, t_i)$ in a set. The difference with DDPG is that the time at each step is also stored. For each control path, we can approximate $\nabla_\theta V^{\pi_\theta}(s_0)$ by computing

$$\sum_{i=0} e^{-\alpha t_i} \{ \nabla_\theta \pi_\theta(s_i) \nabla_a Q(s, a | \omega) |_{a=\pi_\theta(s_i)} + \nabla_\theta e^{-\alpha \tau_\theta(s_i)} Q(s_{i+1}, \pi_\theta(s_{i+1}) | \omega) \} \quad (26)$$

and averaging it over P paths. We can approximate the policy gradient by generating M initial states s_0 from the initial state distribution d_0 and taking the average of this calculation for each of them. If we take P, N and M to be infinitely large, and if $Q(s, a | \omega)$ is a good approximation of $Q^{\pi_\theta}(s, a)$, we can calculate the correct policy gradient.

Algorithm 2 Ideal algorithm for Self-Triggered Control RL

Write me.

4.4 Practical Algorithm

As I mentioned before, the above algorithm does not take into account the problem of computational complexity. Therefore, from now on, we consider an efficient method to approximate the policy gradient. The most important point is the state distribution of the mini-batch which takes the sample mean to approximate equation (24).

Assuming that the update of the actor is very gradual, the replay buffer stores the experience gained by policies similar to the current policy. Thus, for each experience $(s_i, u_i, \tau_i, r_i, t_i)$, if we create a mini-batch E by sampling the experience with probability $e^{-\alpha t_i}$, we can expect that the sample mean

$$\frac{1}{N} \sum_{(s, s') \in E} \{\nabla_{\theta} \pi_{\theta}(s) \nabla_a Q(s, a | \omega)|_{a=\pi_{\theta}(s)} + \nabla_{\theta} e^{-\alpha \tau_{\theta}(s)} Q(s, \pi_{\theta}(s) | \omega)\} \quad (27)$$

for the mini-batch will approximate (24) well. This is because the distribution of the mini-batch is discounted for time t .

Algorithm 3 Practical algorithm for Self-Triggered Control RL

Write me.

for episode = 0 to M **do**

 Initialize s_0

 Set $i = 0, t_i = 0$

while $t_i < T$ **do**

 Select $[u_i, \tau_i] = \pi_{\theta}(s_i) + e_i$

 Execute action u_i for τ_i second to the environment

 Receive r_i and observe s_{i+1}

 Store $(s_i, u_i, \tau_i, r_i, s_{i+1}, t_i)$ to the replay buffer

 Make mini-batch considering probability $e^{-\alpha t_i}$

 Update critic ω to decrease $L = \sum_{(s, u, \tau)}$

 Calculate approximated policy gradient using (27)

 Update actor with approximated policy gradient

 Update target network

end while

end for

We refer to this approach as the proposed method 1.

4.5 Other approach

In DPG for reinforcement learning problems with a fixed discount factor, it is sufficient that the gradient of the critic with respect to a can correctly approximate that of the Q -function. However, in the case of reinforcement learning for self-triggered control considered in this paper, the value of $Q(s, \pi(s))$ itself must also be correctly approximated from the equation (24). Therefore, we need to pay attention to the TD learning

of critic. However, as mentioned in the previous section, it is a big challenge in terms of computational cost to train the critic until it converges at every update of the actor.

The problem of the above problem is that the discount factor is not fixed. In this section, we consider a reinforcement learning method different from the approach taken in the previous section. In other words, we propose a method to find a policy to indirectly reduce the evaluation function.

First, we redefine the value function $V^\pi(s)$ as follows.

$$\begin{aligned} V^\pi(s) &= \sum_{i=0}^{\infty} \gamma^i \int_{t_i}^{t_{i+1}} \mathbb{E}_w[s(t)^\top Qs(t)dt + \tau_i a_i^\top R a_i - \lambda]dt \\ &= \sum_{i=0}^{\infty} \gamma^i \int_0^{\tau_i} \mathbb{E}_w[s(t)^\top Qs(t)dt + \tau_i a_i^\top R a_i - \lambda | s(0) = s_i]dt \end{aligned} \quad (28)$$

In this section, we consider the problem of finding a policy to minimize the expectation of this value function for the state $s \sim d_0$. From the equation (28), if the reward function is given by $r(s, u, \tau) = \int_0^\tau \mathbb{E}_w[s(t)^\top Qs(t)dt - \tau_i u^\top R u + \lambda | s(0) = s]dt$, then we can give a small cost to keep the state near the origin with a long communication interval and a small input energy. However, since the discount factor is fixed, the value function can be reduced by increasing the number of steps with small communication intervals, the parameter λ is used to adjust the value function.

Since the discount factor is fixed, we can simply apply the DDPG to this problem. We refer to this method as the proposed method 2.

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 $0.01(s) \sim 10.0(s)$.

5.1 Linear System

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

$$\dot{s} = As + Bu = \begin{bmatrix} -1 & 4 \\ 2 & -3 \end{bmatrix} s + \begin{bmatrix} 2 \\ 4 \end{bmatrix} u. \quad (29)$$

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

5.1.1 initial policy

For the comparison with the control performance with that of a naively designed model-based self-triggered control law, we use $\pi^{\text{MB}}(s)$ such that

$$\pi^{\text{MB}}(s) = \underset{u, \tau}{\operatorname{argmin}} \left\{ \int_0^\tau s(t)^\top s(t)dt + u^2 - \lambda\tau + V_{\text{LQR}}(s'(s, u, \tau)) \right\} \quad (30)$$

where $s(0) = s$. The evaluation value of π^{MB} is

$$J(\pi^{\text{MB}}) = . \quad (31)$$

In order to use π^{MB} as an initial policy for reinforcement learning, we represent π^{MB} in a neural network by supervised learning of $\pi^{\text{MB}}(s)$ for M randomly generated states s in the state space S .

5.1.2 proposed method 1

5.1.3 proposed method 2

Fig 3 shows a control path with learned policy π_{RL}^l , stating from $s_0 = [3., 3.]$.

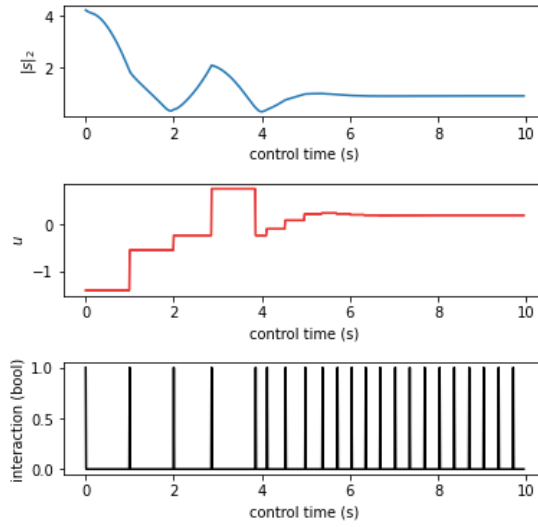


Figure 3: A control path with learned policy π_{RL}^l

Figure 3 shows, from top to bottom, the L2 - norm of state s , 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}^l enlarge the evaluation function $J(\pi)$ compared from the initial policy π_{init}^l . 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

$$J(\pi_{\text{init}}^l) =, J(\pi_{\text{RL}}^l) =, \quad (32)$$

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

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

$$\frac{d}{dt} \begin{pmatrix} \theta \\ \dot{\theta} \end{pmatrix} = \begin{pmatrix} \dot{\theta} \\ \frac{3g}{2l} \sin \theta + \frac{3}{ml^2} a \end{pmatrix}. \quad (33)$$

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

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

5.2.1 Initial Policy

The initial policy π_{init} used in this case is

$$\pi_{\text{init}} = [-Ks \ 0.2] \quad (34)$$

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

5.2.2 proposed method 1

5.2.3 proposed method 2

Fig 4 shows a control path with learned policy, stating from $s_0 = [3., 3.]$.

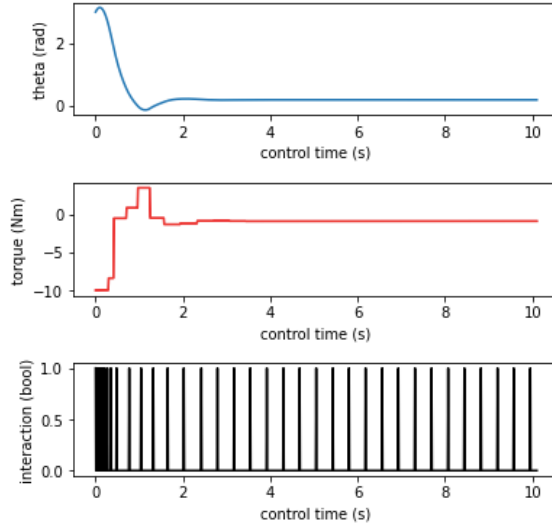


Figure 4: A control path with learned policy π_{RL}

Figure 4 shows, from top to bottom, the angle of the pendulum θ , the control input u ,

and the boolean representing interaction. Similarly to the linear case, Figure 4 shows that the agent decides the communication interval flexibly.

Now let's compare $J(\pi_{\text{RL}})$ and $J(\pi_{\text{init}})$. As in linear case, we approximate $J(\pi) = \mathbb{E}_{s_0}[V^\pi(s_0)]$ in the same way. The result is

$$J(\pi_{\text{init}}) = 4.23 \pm 0.083, \quad J(\pi_{\text{RL}}) = 37.12 \pm 0.071. \quad (35)$$

Since $J(\pi_{\text{init}}) < J(\pi_{\text{RL}})$, we can conclude that the policy has been improved in the non-linear case as well.

6 Conclusion

In this paper, a method for obtaining the state feedback control law is proposed. In particular, it is confirmed by numerical examples that the control law obtained by our method outperforms the control law designed by a naive model-based way.

Furthermore, we also tackle the self-triggered control for nonlinear systems, which has not been considered in previous studies, and successfully obtain useful control laws.

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.
- [8] T. Gommans, D. Antunes, T. Donkers, P. Tabuada, and M. Heemels. “Self-triggered linear quadratic control.” *Automatica*, vol. 50, no. 4, pp. 1279-1287, 2014.

A Appendix

A.1 Model Settings

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 5.

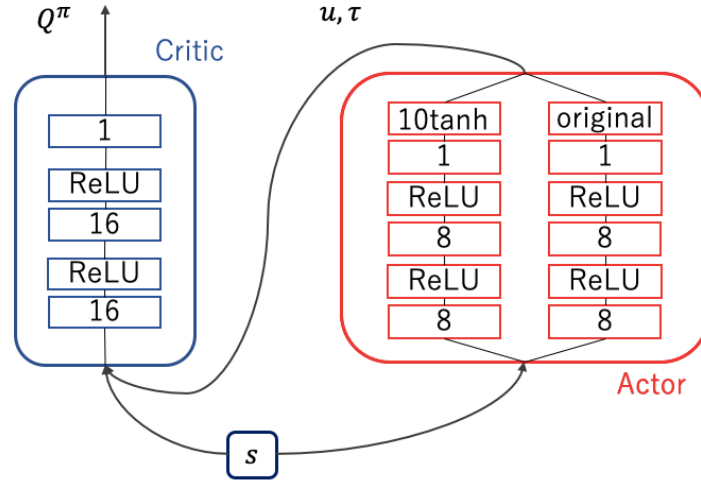


Figure 5: Agent Model

The activation function "original" shown in Fig 5 is defined as $0.99 \times \text{sigmoid} + 0.01$ to meet upper and lower limits of interval described in the next section.

Master's Thesis

Deep Reinforcement Learning for 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 Self-Triggered Control

Ibuki TAKEUCHI

February 2021

Deep Reinforcement Learning for 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.