

## 進捗報告 9.25

## 1 前提

入力アフィン系のセルフトリガー制御を考える.

$$\dot{s} = f(s) + g(s)a \quad (1)$$

## 1.1 倒立振子による実験

倒立時の振子の角度を  $\theta = 0$  とし, 加えられる入力  $A = [-10\text{N} \cdot \text{m}, 10\text{N} \cdot \text{m}]$  と制限されるような倒立振子を考える. この倒立振子のダイナミクスは, 以下のように与えられる.

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

コンピュータで強化学習を行う場合, これを離散化したシステムについて計算を行う必要がある. 上記の状態方程式を離散化すると以下ようになる.

$$\theta_{t+1} = \theta_t + \dot{\theta}_t \delta_t + \frac{3g}{2l} \sin \theta_t \delta_t^2 + \frac{3}{ml^2} a \delta_t^2 \quad (3a)$$

$$\dot{\theta}_{t+1} = \dot{\theta}_t + \frac{3g}{2l} \sin \theta_t \delta_t + \frac{3}{ml^2} a \delta_t \quad (3b)$$

ただし,  $\delta_t$  は離散化定数である.

## 2 現状確認

## 2.1 セルフトリガー制御

図 1 のような制御系を考える.

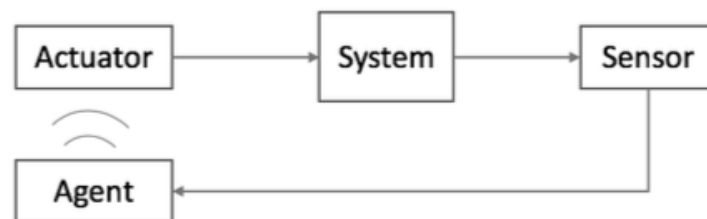


図 1 制御系

これに対するフィードバック制御を考える. 状態変数  $s$  を観測してアクチュエータに入力信号を送信することを「インタラクション」と呼ぶと, セルフトリガー制御では, 連続的なインタラクションは行わずに, 次のイ

インタラク션을何秒後に行うかをエージェントが決定する．それを数式上で表すため，エージェントの制御則  $\pi(s)$  は2つの要素からなるベクトル値関数であるとし，1つ目の要素はアクチュエータに送信する入力  $a$ ，2つ目の要素は次にインタラク션을行うまでの時間間隔  $\tau$  ( $s$ : 秒) を表すものとする．また，次のインタラク션을行う時刻までは1つ前のインタラク션で送信した入力  $a$  を加え続けるものとする (ZOH 制御)．

## 2.2 目標点の確認

研究を通しての目標は「安全性を確保しながら，最適セルフトリガー制御則  $\pi^*$  の強化学習を実現させること」である．ここで

$$\pi^* = \operatorname{argmax}_{\pi} J(\pi) \quad (4)$$

$$J(\pi) = \mathbb{E}_{s_0 \in d_0} [V^{\pi}(s_0)] \quad (5)$$

$$V^{\pi}(s_0) = \sum_{i=0}^{\infty} \gamma^i C_i^{\pi} \quad (6)$$

$$C_i^{\pi} = - \int_{T_i}^{T_{i+1}} s(t)^{\top} Q s(t) dt + \tau_i a_i^{\top} R a_i + \lambda \tau_i, \quad T_i = \sum_{l=0}^i \tau_l \quad (7)$$

であり， $\pi_1, \pi_2$  は  $\pi$  の第1, 第2成分である．また， $i$  はインタラク션の回数を示し， $a_i, \tau_i$  はそれぞれ  $i$  回目のインタラク션での方策  $\pi$  の出力である．

さて，一般的に強化学習では，1ステップ1ステップの行動の良し悪しを評価して方策を更新していく．インタラク션とインタラク션の間の区間を「インターバル」と呼ぶと，式 (6) より，この問題は各インターバルを1ステップとした強化学習問題であると考えることができる．

以下では方策  $\pi$  を  $\theta$  でパラメトライズし， $\theta^* = \operatorname{argmax}_{\theta} J(\pi_{\theta})$  を解くことによって  $\pi^* = \operatorname{argmax}_{\pi} J(\pi)$  を得るものとする．方策勾配を用いた強化学習では  $\nabla_{\theta} J(\pi)$  を用いて  $\theta^*$  を求める．その際方策勾配  $\nabla_{\theta} J(\pi)$  の近似のため，実環境とのインタラク션によって得られたデータ組  $\{s, a, r, s'\}$  を用いる．「学習中の安全」という言葉を，「このデータ組の収集を決められた安全領域  $\mathcal{C}$  の内部でのみ行うこと」と定義する．

## 2.3 実現可能性の検証: サンプル値系での実験

上記の目標を達成する見込みがあるのかを検証するために，サンプル値系での実験を行う．サンプル値系では，セルフトリガー制御と同様に連続的なインタラク션は行わない．セルフトリガー制御との違いは，インタラク션の間隔がエージェントによって状態  $s$  依存で決定するのではなく，制御問題の設定として定数  $t_{\text{int}}$  で与えられる点である．したがってサンプル値系での制御方策  $\pi_{\text{sample}}$  は，アクチュエータに送信する入力信号  $a$  のみを出力する関数として与える．

サンプル値系での実験により， $t_{\text{int}} = 0.001(s)$  のサンプル値系での最適方策

$$\begin{cases} \pi_{\text{sample},1} = - \operatorname{argmax}_{\pi_{\text{sample}}} \sum_{i=0}^{\infty} \int_{it_{\text{int}}}^{(i+1)t_{\text{int}}} s(t)^{\top} Q s(t) dt + t_{\text{int}} a_i^{\top} R a_i \\ a_i = \pi_{\text{sample}}(s(it_{\text{int}})) \end{cases} \quad (8)$$

を初期方策として， $t_{\text{int}} = 0.002(s)$  のサンプル値系での最適方策

$$\begin{cases} \pi_{\text{sample},2} = - \operatorname{argmax}_{\pi_{\text{sample}}} \sum_{i=0}^{\infty} \int_{it_{\text{int}}}^{(i+1)t_{\text{int}}} s(t)^{\top} Q s(t) dt + t_{\text{int}} a_i^{\top} R a_i \\ a_i = \pi_{\text{sample}}(s(it_{\text{int}})) \end{cases} \quad (9)$$

を学習中の安全性を満たしながら学習できるかを検証する．

## 2.4 セルフトリガー制御への発展

前節での検証によって、インタラクション間隔を大きくしても安全強化学習を行うことが可能であることを確認できたとする。サンプル値系での制御則は入力信号  $a$  のみを出力する関数であったので、入力信号  $a$  とインタラクション間隔  $\tau$  の二つの要素を出力する必要があるセルフトリガー制御の初期方策として方策  $\pi_{\text{sample},1}$  をそのまま用いることはできない。

そこで代替策として、

$$\begin{cases} \pi_1(s) = \pi_{\text{sample},1}(s) \\ \pi_2(s) = 0.001 \end{cases} \quad (10)$$

とする方策  $\pi_{\text{init}}$  をセルフトリガー制御の強化学習のための初期方策として用いる。

## 3 安全性の定義

### 3.1 インタラクション間隔 $\tau$ の安全性

ECBF(後から書きます)

### 3.2 入力信号 $a$ の安全性

強化学習ではデータの収集に環境とのインタラクションを行う必要がある。DDPG と呼ばれるアルゴリズムは方策オン型の強化学習とよばれ、データの収集方策に学習中の暫定最適方策を用いる。したがって、学習初期の方策では安全性が保証されないことがしばしばある。この課題を解決するために、制御バリア関数を用いる。

関数  $h(s)$  が以下の条件を満たす時、システム (1) に対する制御バリア関数であるという。

$$\sup_{a \in A} \left\{ \frac{\partial h}{\partial s}(f(s) + g(s)a) + K(h(s)) \right\} \geq 0 \quad (11)$$

ただし、 $K(s)$  はクラス K 関数である。

さて、2.2 節にて登場した安全領域  $\mathcal{C}$  を

$$\mathcal{C} = \{s \in S \mid h(s) \geq 0\} \quad (12)$$

として与える。このとき  $h(s)$  が制御バリア関数であるならば、状態  $s \in \mathcal{C}$  を初期状態とした時、それ以降の全時刻において、状態  $s$  が  $s \in \mathcal{C}$  を満たすようにする入力が存在することを保証する。そのような入力集合は現時刻での状態  $s$  に依存し、

$$U(s) = \left\{ a \in A \mid \frac{\partial h}{\partial s}(f(s) + g(s)a) + K(h(s)) \geq 0 \right\}, \forall s \in \mathcal{C} \quad (13)$$

としてその集合を与える。

学習中の安全性を確保するために、図 2 のようにエージェントの出力を  $U(s)$  の要素に射影するレイヤーを設ける。

当面は、エージェントの出力  $a_\pi$  に最も近い  $U(s)$  の元への射影を考える。

## 4 今後の方針

### 4.1 シミュレーション環境の構築

さて、制御バリア関数による安全性保証は連続システムに対して行われるものである。また、セルフトリガー制御則の第 2 成分  $\tau$  の出力ロジックを勾配によって更新できるように、 $\tau$  を連続値として扱う必要がある。しかし、コンピュータ上でシミュレーションを行うには (1) を時間に関して離散化を行わなくてはならない。

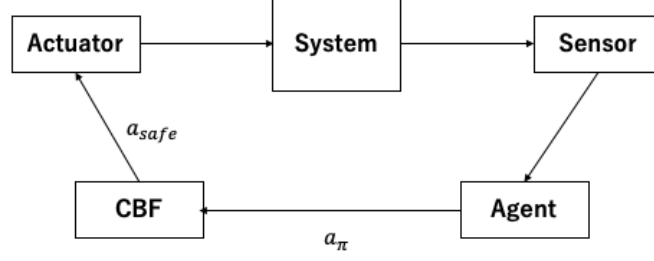


図2 制御系

そこで、インタラクション間隔  $\tau$  を整数個に等間隔に分割し、その時間幅を用いて離散化を行う。この時、離散化幅は  $0.001(s) \sim 0.005(s)$  になるように分割数を調整する。  $\delta_t$  を上から抑える理由は離散化誤差を抑えるためである。  $(\tau)$  を下から抑えるのは、コンピュータ上で  $\frac{1}{\infty} = 0$  になってしまうからで、それを回避するためである。

ここで、式 (7) のインターバル報酬  $C_i^\pi$  が定積分を用いて表されているため、これをシミュレーション環境で近似する手法を考える。ダイナミクス (1) の離散化幅  $\delta_t$  の離散近似システムが

$$s_{t+1} = f_d(s_t, \delta_t) + g_d(s_t, \delta_t)a_t \quad (14)$$

と書かれているとする。インタラクション間隔  $\tau$  を  $N$  分割した時、  $\delta_t = \frac{\tau}{N}$  を用いて

$$C_i^\pi \approx -\delta_t \sum_{k=0}^N s_{n_i+k}^\top Q s_{n_i+k} + \tau_i a_i^\top R a_i + \lambda \tau_i \quad (15)$$

と近似する。ここで  $s_{n_i}$  は  $i$  回目のインタラクションを行った瞬間の状態変数  $s(T_i)$  と同じ値が代入されるものとする。

## 4.2 サンプル値系での安全性確保ロジックの構築

2.3 節で記述した通り、  $t_{\text{int}}$  が  $0.001(s)$  の最適方策  $\pi_{\text{sample},1}$  を初期値として、  $0.002(s)$  の最適方策  $\pi_{\text{sample},2}$  を安全強化学習できるのか検証する。本節ではその安全性の確保方法についてもう少し掘り下げて議論する。

ECBF を用いることによって、入力  $a$  に対するインタラクション間隔  $\tau$  の限界を与えることができる。その値を  $\tau_{\text{max}}(a)$  と書く。もし  $\tau_{\text{max}}(a) < 0.002$  であるなら、CBF を用いて  $U(s)$  の元  $a_{\text{safe}}$  を選び  $\tau_{\text{max}}(a_{\text{safe}}) \geq 0.002$  となれば、次のインタラクションまで  $a_{\text{safe}}$  を加え続けても状態変数が  $\mathcal{C}$  を出ていくことはない。しかし  $\forall a \in U(s)$  に対して  $\tau_{\text{max}}(a) < 0.002$  であるなら、次のインタラクションを  $0.001(s)$  後に行うことで、安全性を確保する必要がある。  $(0.001 \leq \tau_{\text{max}}(a) < 0.002)$  を仮定)

ここまでの議論を整理すると、「サンプル値系における安全性保証」とは「1. 入力信号の安全性、2. サンプル間隔の安全性」を保証することになる。これらが行われることを回避する学習方法については今後検討する。

## 4.3 セルフトリガー制御の強化学習

2.4 節で記した初期方策  $\pi_{\text{init}}$  から、安全性を保証しながら最適セルフトリガー制御則  $\pi^*$  の学習を試みる。サンプル値系とは異なり、インタラクション間隔はエージェントが決定する。したがって、方策関数の出力  $\pi(s) = [a \quad \tau]$  に対して、  $\tau > \tau_{\text{max}}(a)$  となった時には入力信号  $a$  を変更する方法と、通信間隔  $\tau$  を変更する方法の二種類の選択があり、どちらを行うべきか考える必要がある。

## 5 (先行して) セルフトリガー制御の強化学習の実験

現状, ECBF についてまだ深い考察と実装を行っていないので, セルフトリガー制御の安全性の確保は行うことはできない. しかしながら  $\pi_{\text{init}}$  から  $\pi^*$  を強化学習するための環境は整っているため, 安全性については考慮せずに実験を行ってみた.

実験環境は Open-AI Gym の pendulum で, 初期状態が  $\theta \sim N(0, 0.1), \dot{\theta} \sim U(-1, 1)$  となるような環境で行った. また, 10 秒間の制御を 1 エピソードと定義する.

### 5.1 $\lambda$ による学習過程の変化

式 (7) より,  $\tau$  を大きくするモチベーションは  $\lambda$  に大きく依存するため, それを様々な値にすることで学習過程にどのような変化が見られるのかを確認してみた. ここでは,  $\lambda$  を 0.05, 0.5, 5, 50 に設定し, それぞれ 3 回, 200000 ステップ (インタラクションの回数) の強化学習を行った. また,  $\tau$  は 0.001 ~ 0.1 の範囲に制限している.

まずそれぞれの学習タスクの中で, 1 エピソードの間の  $\tau$  の平均値の変化を見てみる.

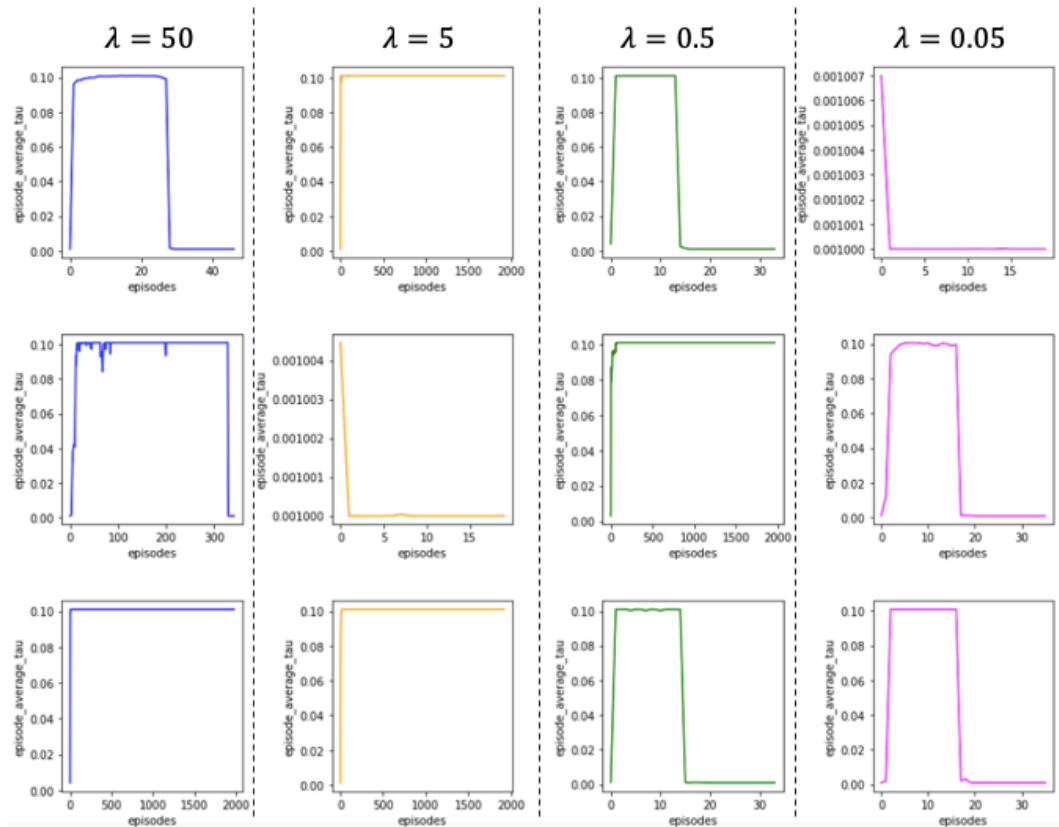


図3 学習を通しての  $\tau$  の変化

ここで注意したいのは,  $\tau$  によって 200000 回のインタラクションに含まれるエピソード数が異なる点である. 例えば  $\tau$  の学習が進まずにずっと  $\tau = 0.001$  であった場合, 10 秒間の制御の間には 10000 回のインタラクションが行われる. したがって 200000 回のインタラクションの間には 20 エピソードしか含まれない. 逆に, 初期から  $\tau$  の学習が進み  $\tau = 0.1$  となることが増えてくると, エピソード数は増えてくる. 図3の  $x$  軸の範囲が異なるのはそのためである.

次に, 制御性能 (倒立振子の制御) の変化を見るためこの時のそれぞれのエピソード報酬の変化を見てみる.

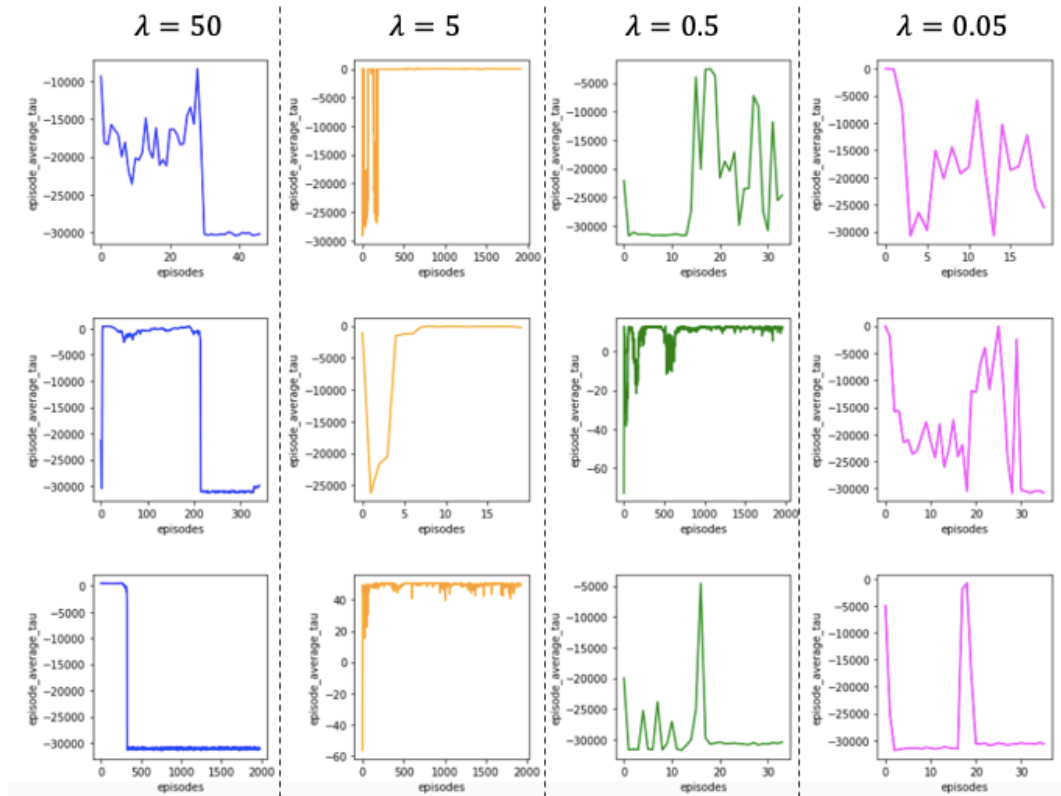


図4 学習を通してのエピソード報酬の変化

## 5.2 考察

図4の各ブロックが表す学習の記録は図3と対応している. これら2つの図を見比べてみるといくつかの点に気付く.

1.  $\lambda = 0.1$  を初期に学習できていないと, 最終的に得られる制御性能自体も悪化する.
2.  $\lambda$  が大きい方が  $\tau$  を大きくする方向に学習が進みやすい (当然) が, 毎回そうなる訳でもない点
3.  $\lambda = 0.1$  を学習できていても,  $\lambda = 0.001$  に逆戻りしてしまうことがある.

1. について, reward の改善が見られないのはエピソード数が少ないためであることは予想がつくが,  $\pi_{init}$  は  $\tau = 0.001$  において, 原点付近で線形化したシステムを安定化する方策であるため, reward が減っていくことは腑に落ちない.

3. については, reward をみると  $\tau = 0.1$  を維持している学習とそうでない学習で,  $\tau = 0.1$  の間の制御性能に差がある点ことがわかり, これが理由だと考える. ただし,  $\lambda = 50$  の学習での逆行は理由がわからない.

## 5.3 学習した方策による制御性能

図4の2行3列目が表す学習によって得られた方策  $\pi_a$  の制御を詳しくみしてみる.

図5を見ると, おおよそ振り子を倒立させ続けられていることがわかる. しかし注目したいのは  $\tau$  の変化である. 本来セルフトリガー制御は状態  $s$  によって様々な値をとるはずであるが, ほぼ一定値となってしまう. このことから「倒立時は長くインターバルを空け, 倒れ出した時に多くインタラクションを行う」というほども高度なことを学習できていない. また, これは果たしてリーズナブルな結果なのか今後明らかにしていきたい.

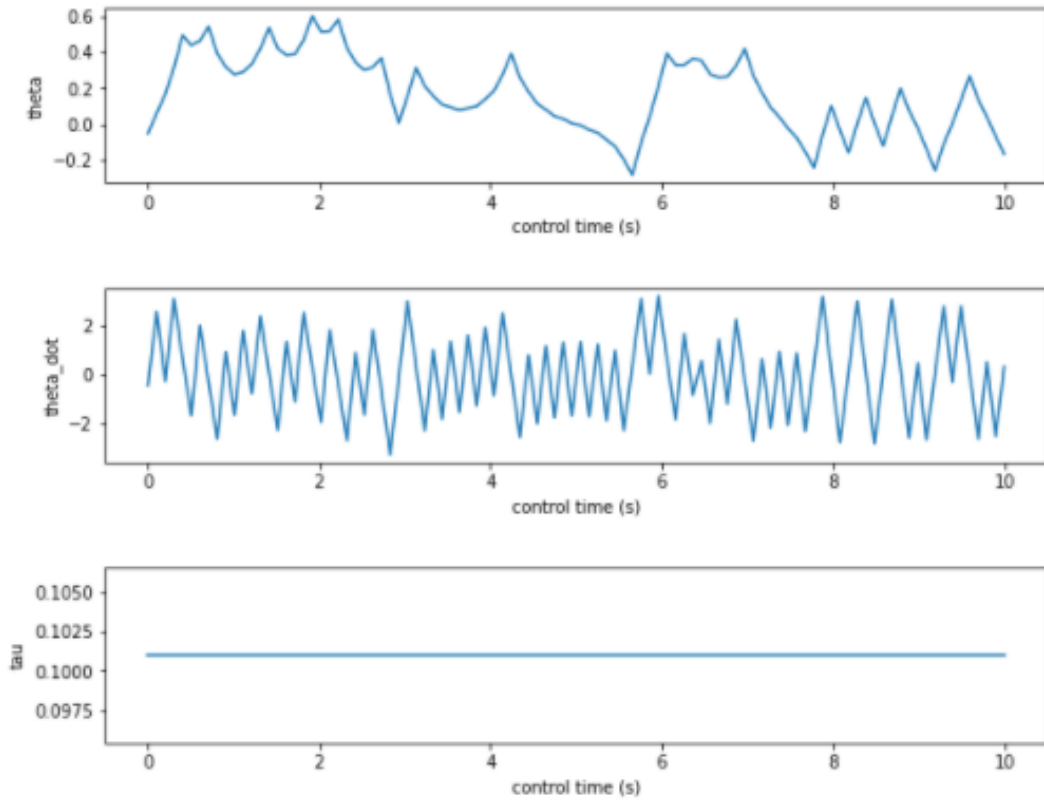


図5  $\pi_a$  による制御ログ

## 6 $\tau(s)$ の学習の様子

前節で、エピソード (10 秒間の制御) 毎にその間の  $\tau(s)$  の平均値を追うと、0.001 から 0.1 に急激に変化していることを確認した。また、5.3 節では、200000 ステップの学習によって得られた方策によって決定される  $\tau(s)$  がインターバルの間ずっと 0.1 となっていることも確認した。

そこで、本節では  $\tau(s)$  を正しく  $s$  依存の関数として学習して、その結果このように「張り付く」ような結果になっているのか否かを確認する。

### 6.1 検証方法

方策関数  $\pi_\theta(s)$  のパラメータ  $\theta$  は 1 ステップごとに行われる。したがってその 1 要素である関数としての  $\tau(s)$  の 1 ステップごとに更新される。ここでは学習途中のとある 100 ステップ分の  $\tau(s)$  を取り出し、 $\theta \in [-\pi, \pi]$ ,  $\dot{\theta} \in [-8, 8]$  の範囲からランダムに 500 個  $s$  を選び、その  $s$  に対する  $\tau(s)$  の変化を確認する。

### 6.2 $\tau(s)$ の変化

まず初めに  $\pi_{\text{init}}$  による  $\tau(s)$  をみしてみる。

これは式 (10) の通りである。

では次に 5.3 節で確認した方策  $\pi_a$  による  $\tau(s)$  をみしてみる。

図 7 をみると、ほとんど全ての状態で  $\pi_a(s) = 0.1$  となっていることがわかる。

では、これらの途中にはどのような変化があったのかを確認する。

図 8 は 6.1 節で述べた、学習中のとある 100 ステップ分の関数  $\tau(s)$  のうち、等間隔に 10 個選んだものをプロットしたものである。図の同じ色の点は同じ関数  $\tau(s)$  による出力を表し、青色から緑色にかけて学習が進んだも

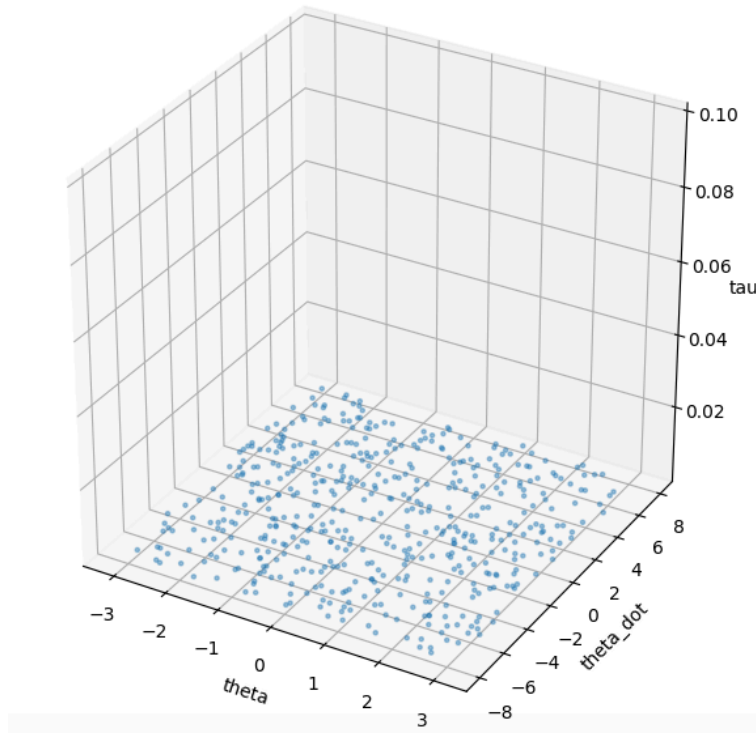


図6  $\pi_{\text{init}}$  の  $\tau(s)$

のを表す. 図8をみると, 少しずつ  $\tau(s) = 0.1$  となる領域が大きくなっていることがわかる.

### 6.3 結果の考察

前節での結果から,  $\tau(s)$  はきちんと  $s$  の関数として学習 (しようと) していることがわかった. では図8で  $\tau(s) = 0.1$  となる領域はどのようにして決まっているのか. それを確認するために, 学習途中にエージェントが経験した状態  $s$  のログを録ってみた. それと  $\tau(s) = 0.1$  となる領域が関係するのか, 図9をみる.

図9でオレンジ色で表しているのは, エージェントが実際に経験した状態変数のパスである (3次元空間で比較しやすいように 'tau' 座標を与えた). これをみると, 経験した状態  $s$  の領域が  $\tau(s) = 0.1$  となる領域となっていることがわかる. これは  $\lambda$  が大きすぎる, 勾配法の学習率が大きすぎる, のどちらかを改善する必要があるそうである.

図3でもみた通り  $\lambda$  を小さくしても  $\tau(s) = 0.1$  に張り付いてしまう傾向は変わらない.

では,  $\tau(s)$  の更新幅を小さくする方法を考える.  $\tau(s)$  は図10のように  $\pi_{\theta}(s)$  を表現するニューラルネットワークの出力の1つとして与えているので,  $\tau(s)$  の学習率のみを小さくすることはできない (入力信号  $a(s)$  の学習も遅くなってしまう).

そこで,  $a(s)$  と  $\tau(s)$  をそれぞれ別のニューラルネットワークで表現してはどうだろうかと考えている.



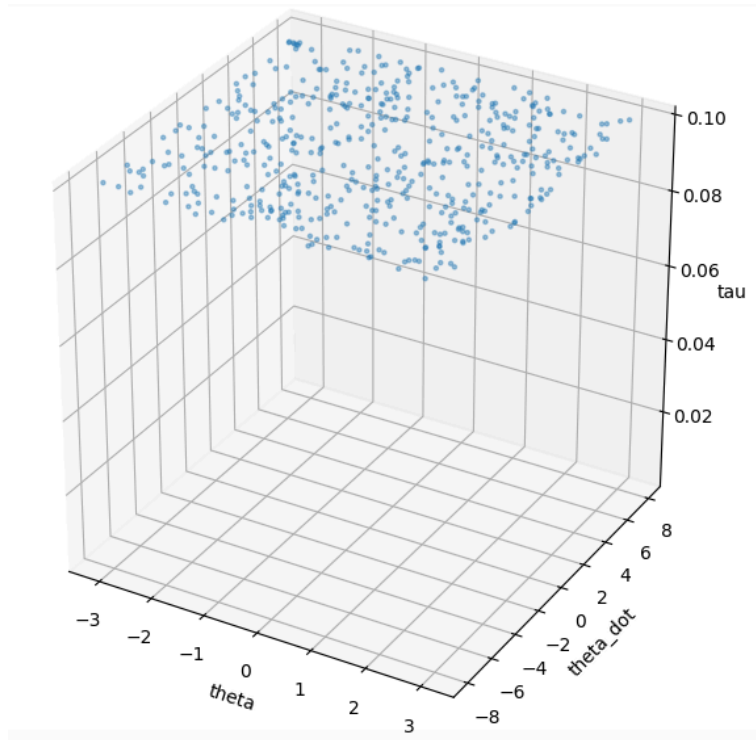


図7  $\pi_a$  の  $\tau(s)$

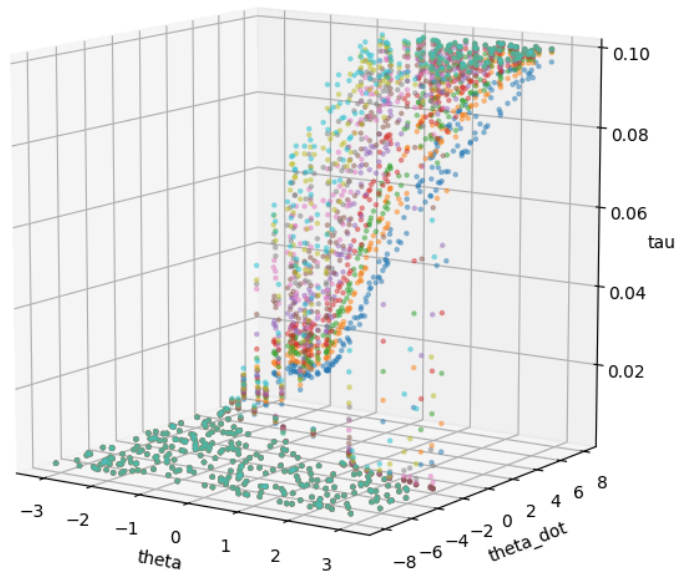


図8 学習が進むにつれてどのように  $\tau(s)$  が変わるのか

## 参考文献

- [1] G. Yang, C. Belta, and R. Tron. “Self-triggered Control for Safety Critical Systems Using Control Barrier Functions.” *In American Control Conference (ACC) Philadelphia, USA*, 2019.

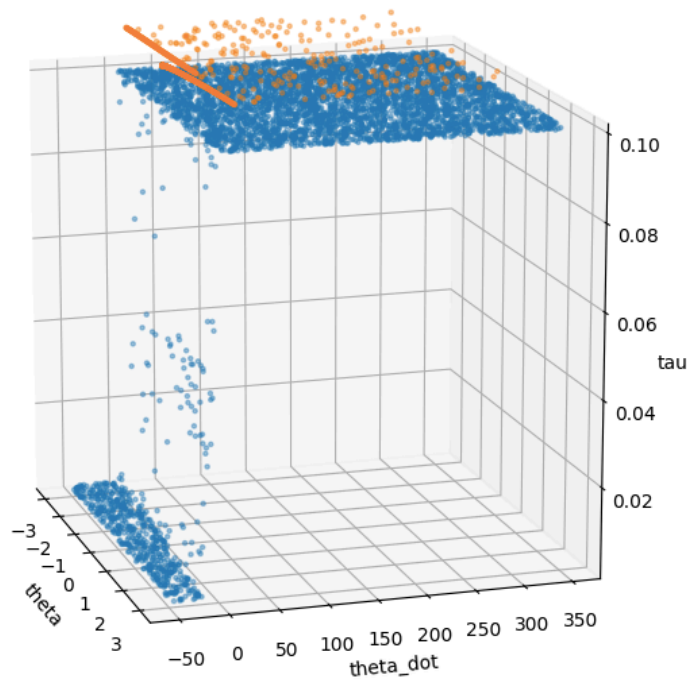


図9 経験した  $s$  と  $\tau(s) = 0.1$  となる領域

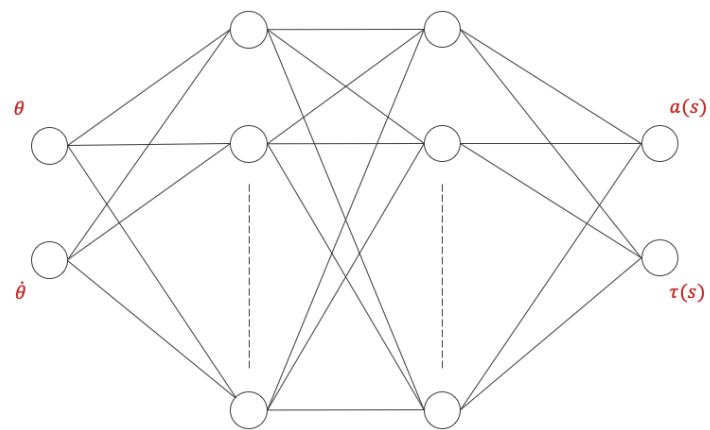


図10 方策  $\pi_{\theta}(s)$

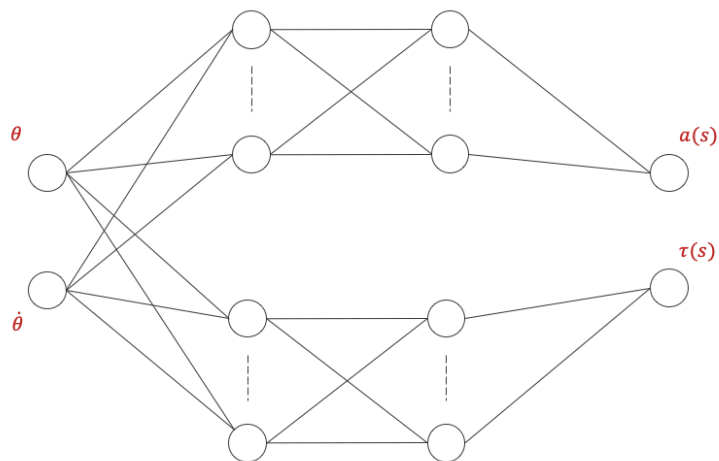


図11 方策  $\pi_{\theta}(s)$  の分離