

دانشگاه صنعتی شریف دانشکدهی مهندسی هوافضا

> پروژه کارشناسی مهندسی کنترل

> > عنوان:

هدایت یادگیری تقویتی مقاوم مبتنی بر بازی دیفرانسیلی در محیطهای پویای چندجسمی با پیشران کم

نگارش:

علی بنی اسد

استاد راهنما:

دكتر هادى نوبهارى

تیر ۱۴۰۱



به نام خدا

دانشگاه صنعتی شریف

دانشکدهی مهندسی هوافضا

پروژه کارشناسی

عنوان: هدایت یادگیری تقویتی مقاوم مبتنی بر بازی دیفرانسیلی در محیطهای پویای چندجسمی با پیشران کم

نگارش: على بنى اسد

كميتهى ممتحنين

استاد راهنما: دكتر هادى نوبهارى امضاء:

استاد مشاور: استاد مشاور

استاد مدعو: استاد ممتحن امضاء:

تاريخ:

سپاس

از استاد بزرگوارم جناب آقای دکتر نوبهاری که با کمکها و راهنماییهای بیدریغشان، بنده را در انجام این پروژه یاری دادهاند، تشکر و قدردانی میکنم. از پدر دلسوزم ممنونم که در انجام این پروژه مرا یاری نمود. در نهایت در کمال تواضع، با تمام وجود بر دستان مادرم بوسه میزنم که اگر حمایت بیدریغش، نگاه مهربانش و دستان گرمش نبود برگ برگ این دست نوشته و پروژه وجود نداشت.

چکیده

در این پژوهش، از یک روش مبتنی بر نظریه بازی به منظور کنترل وضعیت استند سه درجه آزادی چهار پره استفاده شده است. در این روش بازیکن اول سعی در ردگیری ورودی مطلوب می کند و بازیکن دوم با ایجاد اغتشاش سعی در ایجاد خطا در ردگیری بازیکن اول می کند. در این روش انتخاب حرکت با استفاده از تعادل نش که با فرض بدترین حرکت دیگر بازیکن است، انجام می شود. این روش نسبت به اغتشاش ورودی و همچنین نسبت به عدم قطعیت مدل سازی می تواند مقاوم باشد. برای ارزیابی عملکرد این روش ابتدا شبیه سازی هایی در محیط سیمولینک انجام شده است و سپس، با پیاده سازی روی استند سه درجه آزادی صحت عملکرد کنترل کننده تایید شده است.

کلیدواژهها: چهارپره، بازی دیفرانسیلی، نظریه بازی، تعادل نش، استند سه درجه آزادی، مدلمبنا، تنظیمکننده مربعی خطی

¹Game Theory

²Nash Equilibrium

فهرست مطالب

١	یادگیری تقویتی	
١	۱-۱ مفاهيم اوليه	
۲	۱-۱-۱ حالت و مشاهدات	
۲	۲-۱-۱ فضای عمل ۲۰۰۰،۰۰۰،۰۰۰،۰۰۰	
۲	۳-۱-۱ سیاست	
٣	۴-۱-۱ مسیر	
٣	۵-۱-۱ تابع پاداش و بازگشت ،	
۴	۱-۱-۶ ارزش در یادگیری تقویتی	
۵	۲-۱ عامل گرادیان سیاست عمیق قطعی ۲۰۰۰،۰۰۰،۰۰۰،۰۰۰	
۵	۱-۲-۱ یادگیری Q در DDPG	
٧	۲-۲-۱ سیاس <i>ت</i> در DDPG سیاس <i>ت</i> در	
٧	۱-۲-۳ اکتشاف و بهرهبرداری در DDPG	
٨	۴-۲-۱ شبه کد	
١.	۳-۱ عامل گرادیان سیاست عمیق قطعی تاخیری دوگانه	
١.	۱-۴ بهینهسازی سیاست مجاور ۲۰۰۰،۰۰۰،۰۰۰،۰۰۰،۰۰۰	

فهرست جداول

فهرست تصاوير

۱-۱ حلقه تعامل عامل و محیط ۲۰۰۰،۰۰۰ حلقه تعامل عامل و محیط

فصل ۱

يادگيري تقويتي

۱-۱ مفاهیم اولیه

بخشهای اصلی یادگیری تقویتی شامل عامل و محیط است. عامل در محیط قرار دارد و با آن تعامل دارد. در هر مرحله از تعامل بین عامل و محیط، عامل یک مشاهده جزئی از وضعیت محیط انجام می دهد و سپس در مورد اقدامی که باید انجام دهد تصمیم می گیرد. وقتی عامل بر روی محیط عمل می کند، محیط تغییر می کند، اما ممکن است محیط به تنهایی نیز تغییر کند. عامل همچنین یک سیگنال پاداش آز محیط دریافت می کند، عددی که به آن می گوید وضعیت فعلی محیط چقدر خوب یا بد است. هدف عامل به حداکثر رساندن پاداش انباشته خود است که بازگشت نام دارد. یادگیری تقویتی روشهایی هستند که عامل رفتارهای مناسب برای رسیدن به هدف خود را می آموزد. در شکل 1-1 تعامل بین محیط و عامل نشان داده شده است.

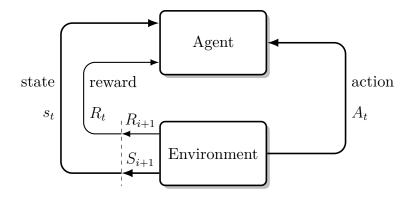
¹Reinforcement Learning (RL)

²Agent

 $^{^3}$ Environment

⁴Reward

⁵Return



شكل ١-١: حلقه تعامل عامل و محيط

۱-۱-۱ حالت و مشاهدات

حالت (s) توصیف کاملی از وضعیت محیط است. همه ی اطلاعات محیط در حالت وجود دارد. مشاهده (o) یک توصیف جزئی از حالت است که ممکن است شامل تمامی اطلاعات نباشد.

۱-۱-۲ فضای عمل

فضای عمل در یادگیری تقویتی، مجموعه ای از تمام اقداماتی است که یک عامل میتواند در محیط خود انجام دهد. این فضا میتواند گسسته $^{\Lambda}$ یا پیوسته $^{\Phi}$ باشد. در این پژوهش فضای عمل پیوسته و در یک بازه مشخص است.

۱-۱-۳ سیاست

یک سیاست^۱ قاعدهای است که یک عامل برای تصمیمگیری در مورد اقدامات خود استفاده میکند. در این پژوهش سیاست قطعی ۱۱ است، که به صورت زیر نشان داده می شود:

$$a_t = \pi(s_t) \tag{1-1}$$

⁶State

⁷Observation

⁸discrete

⁹continuous

¹⁰policy

 $^{^{11}}$ deterministic

در یادگیری تقویتی عمیق از سیاستهای پارامتری شده استفاده می شود. خروجی این سیاستها از توابعی هستند که به مجموعهای از پارامترها (مثلاً وزنها و بایاسهای یک شبکه عصبی) بستگی دارند که می توان آنها را برای تغییر رفتار از طریق برخی الگوریتمهای بهینه سازی تنظیم کرد. در این پژوهش پارامترهای سیاست را با θ نشان داده شده است و سپس نماد آن به عنوان یک زیروند روی سیاست مانند معادله (۲–۱) نشان داده شده است.

$$a_t = \pi_\theta(s_t) \tag{Y-1}$$

١-١-۴ مسير

یک مسیر ۱۲ توالی از حالتها و عملها در محیط است.

$$\tau = (s_0, a_0, s_1, a_1, \cdots) \tag{\Upsilon-1}$$

گذار حالت s+1 به اتفاقاتی که در محیط بین حالت در زمان s و حالت در زمان s+1 میافتد، گفته می شود. این گذارها توسط قوانین طبیعی محیط انجام می شوند و تنها به آخرین اقدام انجام شده توسط عامل (a_t) بستگی دارند. گذار حالت را می توان به صورت زیر تعریف کرد.

$$s_{t+1} = f(s_t, a_t) \tag{Y-1}$$

-1-1 تابع پاداش و بازگشت

تابع پاداش ۱۴ حالت فعلی محیط، آخرین عمل انجام شده و حالت بعدی محیط بستگی دارد. تابع پاداش را میتوان بهصورت زیر تعریف کرد.

$$r_t = R(s_t, a_t, s_{t+1}) \tag{Δ-1}$$

در این پژوهش پاداش تنها تابعی از جفت حالت-عمل ($r_t = R(s_t, a_t)$) است. هدف عامل این است که مجموع پاداشهای به دست آمده در طول یک مسیر را به حداکثر برساند، اما این مفهوم می تواند چند معنی داشته باشد. در این پژوهش این موارد را با نماد $R(\tau)$ نشان داده شده است و به آن تابع بازگشت گفته می شود. یکی از انواع بازگشت، بازگشت بدون تنزیل با افق محدود $R(\tau)$ است که مجموع پاداشهای به دست آمده در یک

¹²Trajectory

¹³state transition

¹⁴reward function

¹⁵Return

¹⁶Finite-Horizon Undiscounted Return

بازه زمانی ثابت از مسیر بهصورت زیر است.

$$R(\tau) = \sum_{t=0}^{T} r_t \tag{9-1}$$

نوع دیگری از بازگشت، بازگشت تنزیل شده با افق نامحدود ۱۷ است که مجموع همه پاداشهایی است که تا به حال توسط عامل به دست آمده است، اما با در نظر گرفتن فاصله زمانی ای که تا دریافت آن پاداش وجود داشته، تنزیل ۱۸ شده است. این فرمول پاداش شامل یک فاکتور تنزیل ۱۹ با نماد γ است.

$$R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t \tag{V-1}$$

۱-۱-۶ ارزش در یادگیری تقویتی

در یادگیری تقویتی، دانستن ارزش^۲ یک حالت یا جفت حالت عمل ضروری است. منظور از ارزش، بازگشت مورد انتظار^{۲۱} است، یعنی اگر از آن حالت یا جفت حالت عمل شروع شود و سپس برای همیشه طبق یک سیاست خاص عمل شود، به طور میانگین چه مقدار پاداش دریافت خواهد کرد. توابع ارزش به شکلی در تقریبا تمام الگوریتمهای یادگیری تقویتی به کار میروند. در اینجا به چهار تابع مهم اشاره میکنیم.

۱. تابع ارزش تحت سیاست $(V^{\pi}(s))$: این تابع، بازگشت مورد انتظار را در صورتی که از حالت s شروع شود و همیشه طبق سیاست π عمل شود، خروجی می دهد.

$$V^{\pi}(s) = \underset{\tau \sim \pi}{\mathbb{E}} [R(\tau)|s_0 = s] \tag{A-1}$$

۱۰ تابع ارزش—عمل تحت سیاست $(Q^{\pi}(s,a))$: این تابع، بازگشت مورد انتظار را در صورتی که از حالت s شروع شود، یک اقدام دلخواه a (که ممکن است از سیاست π نباشد) انجام شود و سپس برای همیشه طبق سیاست π عمل شود، خروجی می دهد.

$$Q^{\pi}(s,a) = \mathbb{E}_{\tau \sim \pi}[R(\tau)|s_0 = s, a_0 = a]$$
 (9-1)

۳. تابع ارزش بهینه $(V^*(s))$: این تابع، بازگشت مورد انتظار را در صورتی که از حالت s شروع شود و همیشه طبق سیاست بهینه در محیط عمل شود، خروجی میدهد.

 $^{^{17} {\}rm Infinite\text{-}Horizon}$ Discounted Return

¹⁸Discount

¹⁹Discount Factor

²⁰Value

²¹Expected Return

²²On-Policy Value Function

²³On-Policy Action-Value Function

²⁴Optimal Value Function

$$V^*(s) = \max(V^{\pi}(s)) \tag{1.9-1}$$

۴. تابع ارزش—عمل بهینه $(Q^*(s,a))^{(1)}$: این تابع، بازگشت مورد انتظار را در صورتی که از حالت a شروع شود، یک اقدام دلخواه a انجام شود و سپس برای همیشه طبق سیاست بهینه در محیط عمل شود، خروجی می دهد.

$$Q^*(s, a) = \max_{\pi}(Q^{\pi}(s, a))$$
 (11-1)

۲-۱ عامل گرادیان سیاست عمیق قطعی

گرادیان سیاست عمیق قطعی 77 الگوریتمی است که همزمان یک تابع Q و یک سیاست را یاد می گیرد. این الگوریتم برای الگوریتم برای یادگیری تابع Q از داده های غیرسیاست محور 77 و معادله بلمن استفاده می کند. این الگوریتم برای یادگیری سیاست نیز از تابع Q استفاده می کند.

این رویکرد وابستگی نزدیکی به یادگیری Q دارد. اگر تابع ارزش-عمل بهینه مشخص باشد، در هر حالت داده شده عمل بهینه را می توان با حل کردن معادله (1-1) به دست آورد.

$$a^*(s) = \arg\max_{a} Q^*(s, a) \tag{17-1}$$

الگوریتم DDPG ترکیبی از یادگیری تقریبی برای $Q^*(s,a)$ و یادگیری تقریبی برای $a^*(s)$ است و به نحوی طراحی شده است که برای محیطهایی با فضاهای عمل پیوسته مناسب باشد. روش محاسبه $a^*(s)$ در این الگوریتم آن را برای فضای پیوسته مناسب میکند. از آنجا که فضای عمل پیوسته است، فرض می شود که تابع الگوریتم آن را برای فضای پیوسته مناسب میکند. از آنجا که فضای عمل پیوسته است، فرض می شود که تابع $Q^*(s,a)$ نسبت به آرگومان عمل مشتق پذیر است. مشتق پذیری این امکان را می دهد که یک روش یادگیری مبتنی بر گرادیان برای سیاست u(s) استفاده شود. سپس، به جای اجرای یک بهینه سازی زمان بر در هر بار محاسبه u(s) می توان آن را با رابطه u(s) را با رابطه u(s) ستقریب زد.

۱-۲-۱ یادگیری Q در DDPG

معادله بلمن که تابع ارزش عمل بهینه $(Q^*(s,a))$ را توصیف میکند، در پایین آورده شدهاست.

$$Q^*(s,a) = \mathop{\mathbf{E}}_{s' \sim P} \left[r(s,a) + \gamma \max_{a'} Q^*(s',a') \right] \tag{17-1}$$

 $^{^{25}}$ Optimal Action-Value Function

²⁶Deep Deterministic Policy Gradient (DDPG)

²⁷Off-Policy

جمله $P(\cdot|s,a)$ به این معنی است که وضعیت بعدی s' توسط محیط از توزیع احتمال $P(\cdot|s,a)$ نمونه گرفته می شود. معادله بلمن نقطه شروع برای یادگیری $Q^*(s,a)$ با یک مقداردهی تقریبی است. پارامترهای یک شبکه عصبی $Q_{\phi}(s,a)$ با علامت ϕ نشان داده شده است. یک مجموعه D از تغییر از یک حالت به حالت دیگر شبکه عصبی $Q_{\phi}(s,a)$ با علامت $Q_{\phi}(s,a)$ نشان می دهد که آیا وضعیت $S_{\phi}(s,a)$ بایانی است یا خیر) جمع آوری شده است. یک تابع خطای میانگین مربعات بلمن $S_{\phi}(s,a)$ استفاده شده است که معیاری برای نزدیکی $S_{\phi}(s,a)$ برای برآورده کردن معادله بلمن است.

$$L(\phi, \mathcal{D}) = \mathop{\mathbf{E}}_{(s, a, r, s', d) \sim \mathcal{D}} \left[\left(Q_{\phi}(s, a) - \left(r + \gamma (1 - d) \max_{a'} Q_{\phi}(s', a') \right) \right)^{2} \right]$$
 (14-1)

در الگوریتم DDPG دو ترفند برای عمکرد بهتر استفاده شدهاست که در ادامه به بررسی آن پرداخته شدهاست.

بافرهای بازی

الگوریتمهای یادگیری تقویتی جهت آموزش یک شبکه عصبی عمیق برای تقریب $Q^*(s,a)$ از بافرهای بازی $^{7\Lambda}$ تجربه شده استفاده می کنند. این مجموعه \mathcal{D} شامل تجربیات قبلی است. برای داشتن رفتار پایدار در الگوریتم، بافر بازی باید به اندازه کافی بزرگ باشد تا شامل یک دامنه گسترده از تجربیات شود. انتخاب داده های بافر به دقت انجام شده است چرا که اگر فقط از داده های بسیار جدید استفاده شود، بیش برازش $^{7\Lambda}$ رخ می دهید و اگر از تجربه بیش از حد استفاده شود، ممکن است فرآیند یادگیری کند شود.

• شبکههای هدف

الگوریتمهای یادگیری Q از شبکههای هدف استفاده میکنند. اصطلاح زیر به عنوان هدف شناخته می شود.

$$r + \gamma(1 - d) \max_{a'} Q_{\phi}(s', a') \tag{12-1}$$

در هنگام کمینه کردن تابع خطای میانگین مربعات بلمن، سعی شده است تا تابع $\mathbb Q$ شبیه تر به این هدف یعنی رابطه (1-1) شود. اما مشکل این است که هدف بستگی به پارامترهای در حال آموزش ϕ دارد. این باعث ایجاد ناپایداری در کمینه کردن تابع خطای میانگین مربعات بلمن می شود. راه حل آن استفاده از یک مجموعه پارامترهایی که با تأخیر زمانی به ϕ نزدیک می شوند. به عبارت دیگر، یک شبکه دوم ایجاد می شود که به آن شبکه هدف گفته می شود. شبکه هدف دنباله ی شبکه اول را دنبال می کند. پارامترهای شبکه هدف با نشان ϕ_{targ} نشان داده می شوند. در الگوریتم DDPG، شبکه هدف در هر

²⁸Replay Buffers

²⁹Overfit

بهروزرسانی شبکه اصلی، با میانگینگیری پولیاک^{۳۰} بهروزرسانی میشود.

$$\phi_{\text{targ}} \leftarrow \rho \phi_{\text{targ}} + (1 - \rho)\phi \tag{19-1}$$

در رابطه بالا ϕ یک فراپارامتر $^{"}$ است که بین صفر و یک انتخاب می شود. در این پژوهش این مقدار نزدیک به یک درنظرگرفته شدهاست.

الگوریتم DDPG نیاز به یک شبکه سیاست هدف $(\mu_{\theta_{targ}})$ برای محاسبه عملهایی که به طور تقریبی بیشینه DDPG نیاز به یک شبکه سیاست هدف از همان روشی که تابع Q به دست می آید یعنی با میانگین گیری پولیاک از پارامترهای سیاست در طول زمان آموزش استفاده می شود.

با درنظرگرفتن موارد اشارهشده، یادگیری Q در DDPG با کمینه کردن تابع خطای میانگین مربعات بلمن (MSBE) یعنی معادله (۱۷–۱) با استفاده از کاهش گرادیان تصادفی 77 انجام میشود.

$$L(\phi, \mathcal{D}) = \mathop{\mathbf{E}}_{(s, a, r, s', d) \sim \mathcal{D}} \left[\left(Q_{\phi}(s, a) - \left(r + \gamma (1 - d) Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s')) \right) \right)^{2} \right]$$
 (1V-1)

۲-۲-۱ ساست در DDPG

در این بخش یک سیاست تعیینشده $\mu_{\theta}(s)$ یادگرفته می شود تا عملی را انجام می دهد که بیشینه $Q_{\phi}(s,a)$ رخ دهد. از آنجا که فضای عمل پیوسته است و فرض شده است که تابع Q نسبت به عمل مشتق پذیر است، معادله زیر با استفاده از صعود گرادیان q (تنها نسبت به پارامترهای سیاست) حل می شود.

$$\max_{\theta} \mathop{\mathbf{E}}_{s \sim \mathcal{D}} \left[Q_{\phi}(s, \mu_{\theta}(s)) \right] \tag{1A-1}$$

۱-۲-۳ اکتشاف و بهرهبرداری در DDPG

برای بهبود اکتشاف^{۳۴} در سیاستهای DDPG، در زمان آموزش نویز به عملها اضافه میشود. نویسندگان مقاله اصلی DDPG توصیه کردهاند که نویز ^{۳۵}OU با زمانبندی همارتباطی^{۳۶} اضافه شود. در زمان سنجش بهرهبرداری^{۳۷} سیاست از آنچه یادگرفته است، نویز به عملها اضافه نمیشود.

 $^{^{30}}$ Polyak Averaging

 $^{^{31}}$ Hyperparameter

³²Stochastic Gradient Descent

³³Gradient Ascent

 $^{^{34}}$ Exploration

³⁵Ornstein–Uhlenbeck

 $^{^{36}\}mathrm{Time\text{-}Correlated}$

 $^{^{37}}$ Exploitation

۱-۲-۱ شبه کد

در این بخش الگوریتم DDPG پیادهسازی شده آورده شده است. در این پژوهش الگوریتم ۱ در محیط پایتون با استفاده از کتابخانه TensorFlow پیادهسازی شدهاست.

 (\mathcal{D}) ورودی: پارامترهای اولیه سیاست (θ) ، پارامترهای تابع (ϕ) ، بافر بازی خالی

 $\phi_{\mathrm{targ}} \leftarrow \phi$ ، $\theta_{\mathrm{targ}} \leftarrow \theta$ دهید قرار دهید ایرامترهای پارامترهای هدف را برابر با

تا وقتی همگرایی رخ دهد:

وضعیت $a=\mathrm{clip}(\mu_{\theta}(s)+\epsilon,a_{\mathrm{Low}},a_{\mathrm{High}})$ وضعیت $a=\mathrm{clip}(\mu_{\theta}(s)+\epsilon,a_{\mathrm{Low}},a_{\mathrm{High}})$ و عمل $\epsilon\sim\mathcal{N}$

عمل a را در محیط اجرا کنید.

وضعیت بعدی s'، پاداش r و سیگنال پایان d را مشاهده کنید تا نشان دهد آیا s' پایانی است یا خیر. اگر s' پایانی است، وضعیت محیط را بازنشانی کنید.

اگر زمان بهروزرسانی فرا رسیده است:

به ازای هر تعداد بهروزرسانی:

یک دسته تصادفی گذر از یک حالت به حالت دیگر، $B = \{(s,a,r,s',d)\}$ ، از \mathcal{D} نمونهگیری شود.

اهداف را محاسبه كنيد:

$$y(r, s', d) = r + \gamma (1 - d) Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

تابع Q را با یک مرحله از نزول گرادیان با استفاده از رابطه زیر بهروزرسانی کنید:

$$\nabla_{\phi} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi}(s,a) - y(r,s',d))^2$$

سیاست را با یک مرحله از صعود گرادیان با استفاده از رابطه زیر بهروزرسانی کنید:

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi}(s, \mu_{\theta}(s))$$

شبکههای هدف را با استفاده از معادلات زیر بهروزرسانی کنید:

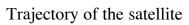
$$\phi_{\text{targ}} \leftarrow \rho \phi_{\text{targ}} + (1 - \rho)\phi$$

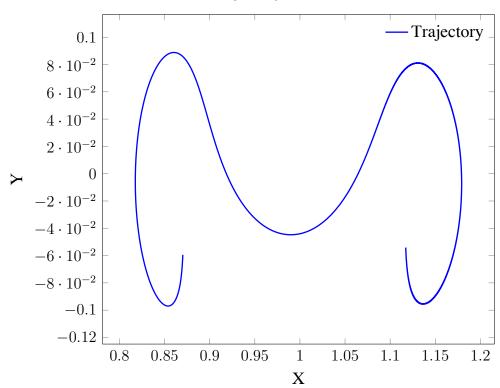
$$\theta_{\text{targ}} \leftarrow \rho \theta_{\text{targ}} + (1 - \rho)\theta$$

۱-۳ عامل گرادیان سیاست عمیق قطعی تاخیری دوگانه

۱-۴ بهینهسازی سیاست مجاور

۱-۵ سندباکس





Bibliography

Abstract

In this study, a quadcopter stand with three degrees of freedom was controlled using game theory-based control. The first player tracks a desired input, and the second player creates a disturbance in the tracking of the first player to cause an error in the tracking. The move is chosen using the Nash equilibrium, which presupposes that the other player made the worst move. In addition to being resistant to input interruptions, this method may also be resilient to modeling system uncertainty. This method evaluated the performance through simulation in the Simulink environment and implementation on a three-degree-of-freedom stand.

Keywords: Quadcopter, Differential Game, Game Theory, Nash Equilibrium, Three Degree of Freedom Stand, Model Base Design, Linear Quadratic Regulator



Sharif University of Technology Department of Aerospace Engineering

Bachelor Thesis

Robust Reinforcement Learning Differential Game Guidance in Low-Thrust, Multi-Body Dynamical Environments

By:

Ali BaniAsad

Supervisor:

Dr.Hadi Nobahari

July 2022