



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:

دکتر هادی نوبهاری

تیر ۱۴۰۱

سلام

## سپاس

از استاد بزرگوالم 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<sup>۱</sup> به منظور کنترل وضعیت استند سه درجه آزادی چهارپره استفاده شده است. در این روش بازیکن اول سعی در ردگیری ورودی مطلوب می‌کند و بازیکن دوم با ایجاد اغتشاش سعی در ایجاد خطا در ردگیری بازیکن اول می‌کند. در این روش انتخاب حرکت با استفاده از تعادل نش<sup>۲</sup> که با فرض بدترین حرکت دیگر بازیکن است، انجام می‌شود. این روش نسبت به اغتشاش ورودی و همچنین نسبت به عدم قطعیت مدل‌سازی می‌تواند مقاوم باشد. برای ارزیابی عملکرد این روش ابتدا شبیه‌سازی‌هایی در محیط سیمولینک انجام شده است و سپس، با پیاده‌سازی روی استند سه درجه آزادی صحت عملکرد کنترل‌کننده تایید شده است.

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

---

<sup>1</sup>Game Theory

<sup>2</sup>Nash Equilibrium

# فهرست مطالب

۲	۱ یادگیری تقویتی
۲	۱-۱ مفاهیم اولیه
۳	۱-۱-۱ حالت و مشاهدات
۳	۱-۱-۲ فضای عمل
۴	۱-۱-۳ سیاست
۴	۱-۱-۴ مسیر
۵	۱-۱-۵ تابع پاداش و بازگشت
۵	۱-۱-۶ ارزش در یادگیری تقویتی
۶	۲-۱ عامل گرادیان سیاست عمیق قطعی
۷	۳-۱ عامل TD3

# فهرست جداول

# فهرست تصاویر

۱-۱ حلقه تعامل عامل و محیط	۳
----------------------------	---

# فصل ۱

## یادگیری تقویتی

### ۱-۱ مفاهیم اولیه

بخش‌های اصلی یادگیری تقویتی<sup>۱</sup> شامل عامل<sup>۲</sup> و محیط<sup>۳</sup> 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 می‌دهد و سپس در مورد اقدامی که باید انجام دهد تصمیم می‌گیرد. وقتی عامل بر روی محیط عمل می‌کند، محیط تغییر می‌کند، اما ممکن است محیط به تنهایی نیز تغییر کند. عامل همچنین یک سیگنال پاداش<sup>۴</sup> از محیط دریافت می‌کند، عددی که به آن می‌گویند وضعیت فعلی محیط چقدر خوب یا بد است. هدف عامل به حداکثر رساندن پاداش انباشته خود است که بازگشت<sup>۵</sup> نام دارد. یادگیری تقویتی روش‌هایی هستند که عامل رفتارهای مناسب برای رسیدن به هدف خود را می‌آموزد. در شکل ۱-۱ تعامل بین محیط و عامل نشان داده شده است.

---

<sup>1</sup>Reinforcement Learning (RL)

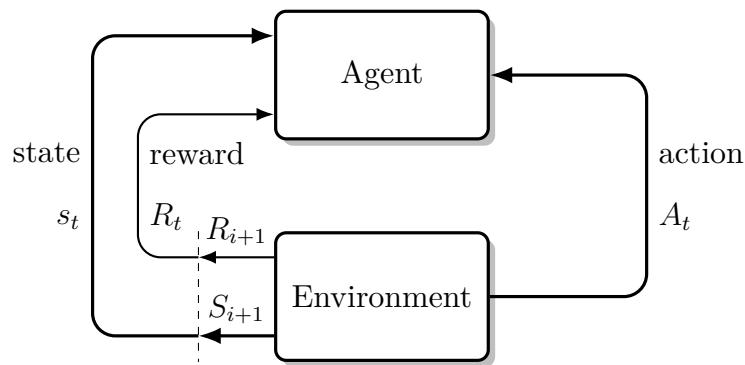
<sup>2</sup>Agent

<sup>3</sup>Environment

<sup>4</sup>Reward

<sup>5</sup>Return





شکل ۱-۱: حلقه تعامل عامل و محیط

### ۱-۱-۱ حالت و مشاهدات

حالت<sup>۶</sup> ( $s$ ) توصیف کاملی از وضعیت محیط است. همه‌ی اطلاعات محیط در حالت وجود دارد. مشاهده<sup>۷</sup> ( $o$ ) یک توصیف جزئی از حالت است که ممکن است تمامی اطلاعات نباشد.

### ۲-۱-۱ فضای عمل

فضای عمل در یادگیری تقویتی، مجموعه‌ای از تمام اقداماتی است که یک عامل می‌تواند در محیط خود انجام دهد. این فضا می‌تواند گسسته<sup>۸</sup> یا پیوسته<sup>۹</sup> باشد. در این پژوهش فضای عمل پیوسته و در یک بازه مشخص است.

<sup>۶</sup>State

<sup>۷</sup>Observation

<sup>۸</sup>discrete

<sup>۹</sup>continuous

## ۳-۱-۱ سیاست

یک سیاست<sup>۱۰</sup> قاعده‌ای است که یک عامل برای تصمیم‌گیری در مورد اقدامات خود استفاده می‌کند. در این پژوهش سیاست قطعی<sup>۱۱</sup> است، که به صورت زیر نشان داده می‌شود:

$$a_t = \pi(s_t) \quad (۱-۱)$$

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

$$a_t = \pi_\theta(s_t) \quad (۲-۱)$$

## ۴-۱-۱ مسیر

یک مسیر<sup>۱۲</sup> توالی‌ای از حالت‌ها و عمل‌ها در محیط است.

$$\tau = (s_0, a_0, s_1, a_1, \dots) \quad (۳-۱)$$

گذار حالت<sup>۱۳</sup> به اتفاقاتی که در محیط بین حالت در زمان  $s$  و حالت در زمان  $s + 1$  می‌افتد، گفته می‌شود. این گذارها توسط قوانین طبیعی محیط انجام می‌شوند و تنها به آخرین اقدام انجام شده توسط عامل ( $a_t$ ) بستگی دارند. گذار حالت را می‌توان به صورت زیر تعریف کرد.

$$s_{t+1} = f(s_t, a_t) \quad (۴-۱)$$

<sup>10</sup>policy<sup>11</sup>deterministic<sup>12</sup>Trajectory<sup>13</sup>state transition

## ۵-۱-۱ تابع پاداش و بازگشت

تابع پاداش<sup>۱۴</sup> حالت فعلی محیط، آخرین عمل انجام شده و حالت بعدی محیط بستگی دارد. تابع پاداش را می‌توان به صورت زیر تعریف کرد.

$$r_t = R(s_t, a_t, s_{t+1}) \quad (5-1)$$

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 می‌تواند چند معنی داشته باشد. در این پژوهش این موارد را با نماد  $R(\tau)$  نشان داده شده است و به آن تابع بازگشت<sup>۱۵</sup> گفته می‌شود. یکی از انواع بازگشت، بازگشت بدون تنزیل با افق محدود<sup>۱۶</sup> است که مجموع پاداش‌های به دست آمده در یک بازه زمانی ثابت از مسیر به صورت زیر است.

$$R(\tau) = \sum_{t=0}^T r_t \quad (6-1)$$

نوع دیگری از بازگشت، بازگشت تنزیل شده با افق نامحدود<sup>۱۷</sup> است که مجموع همه پاداش‌هایی است که تا به حال توسط عامل به دست آمده است، اما با در نظر گرفتن فاصله زمانی‌ای که تا دریافت آن پاداش وجود داشته، تنزیل<sup>۱۸</sup> شده است. این فرمول پاداش شامل یک فاکتور تنزیل<sup>۱۹</sup> با نماد  $\gamma$  است.

$$R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t \quad (7-1)$$

## ۶-۱-۱ ارزش در یادگیری تقویتی

در یادگیری تقویتی، دانستن ارزش<sup>۲۰</sup> یک حالت یا جفت حالت-عمل ضروری است. منظور از ارزش، بازگشت مورد انتظار<sup>۲۱</sup> 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 تقریباً تمام الگوریتم‌های یادگیری تقویتی به کار می‌روند. در اینجا به چهار تابع مهم اشاره می‌کنیم.

<sup>14</sup>reward function

<sup>15</sup>Return

<sup>16</sup>Finite-Horizon Undiscounted Return

<sup>17</sup>Infinite-Horizon Discounted Return

<sup>18</sup>Discount

<sup>19</sup>Discount Factor

<sup>20</sup>Value

<sup>21</sup>Expected Return

۱. تابع ارزش تحت سیاست<sup>۲۲</sup> ( $V^\pi(s)$ ): این تابع، بازگشت مورد انتظار را در صورتی که از حالت  $s$  شروع شود و همیشه طبق سیاست  $\pi$  عمل شود، خروجی می‌دهد.

$$V^\pi(s) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s] \quad (۸-۱)$$

۲. تابع ارزش-عمل تحت سیاست<sup>۲۳</sup> ( $Q^\pi(s, a)$ ): این تابع، بازگشت مورد انتظار را در صورتی که از حالت  $s$  شروع شود، یک اقدام دلخواه  $a$  (که ممکن است از سیاست  $\pi$  نباشد) انجام شود و سپس برای همیشه طبق سیاست  $\pi$  عمل شود، خروجی می‌دهد.

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi} [R(\tau) | s_0 = s, a_0 = a] \quad (۹-۱)$$

۳. تابع ارزش بهینه<sup>۲۴</sup> ( $V^*(s)$ ): 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.

$$V^*(s) = \max_{\pi} (V^\pi(s)) \quad (۱۰-۱)$$

۴. تابع ارزش-عمل بهینه<sup>۲۵</sup> ( $Q^*(s, a)$ ): این تابع، بازگشت مورد انتظار را در صورتی که از حالت  $s$  شروع شود، یک اقدام دلخواه  $a$  انجام شود و سپس برای همیشه طبق سیاست بهینه در محیط عمل شود، خروجی می‌دهد.

$$Q^*(s, a) = \max_{\pi} (Q^\pi(s, a)) \quad (۱۱-۱)$$

## ۲-۱ عامل گرادیان سیاست عمیق قطعی

گرادیان سیاست عمیق قطعی<sup>۲۶</sup> الگوریتمی است که همزمان یک تابع  $Q$  و یک سیاست را یاد می‌گیرد. این الگوریتم برای یادگیری تابع  $Q$  از داده‌های غیرسیاست محور<sup>۲۷</sup> و معادله بلمن استفاده می‌کند. این الگوریتم برای یادگیری سیاست نیز از تابع  $Q$  استفاده می‌کند.

<sup>22</sup>On-Policy Value Function

<sup>23</sup>On-Policy Action-Value Function

<sup>24</sup>Optimal Value Function

<sup>25</sup>Optimal Action-Value Function

<sup>26</sup>Deep Deterministic Policy Gradient (DDPG)

<sup>27</sup>Off-Policy

این رویکرد وابستگی نزدیکی به یادگیری  $Q$  دارد. اگر تابع ارزش-عمل بهینه را مشخص باشد، در هر حالت داده شده، عمل بهینه را می‌توان با حل کردن معادله (۱۲-۱) به دست آورد.

$$a^*(s) = \arg \max_a Q^*(s, a) \quad (۱۲-۱)$$

### ۳-۱ عامل TD3

# 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