



Available in:  
Journal.isrc.ac.ir

Journal of  
Space Science, Technology  
& Applications (Persian)

Vol. 4, No. 1, pp.: 64-77  
2024

DOI:  
10.22034/jssta.2024.446359.1153

### Article Info

Received: 2024-03-03

Accepted: 2024-05-18

### Keywords

Low-Thrust , Equinoctial  
orbital elements ,  
Reinforcement Learning ,  
Actor, Critic networks ,  
Agent

### How to Cite this article

Hamed. Soleymani, M. Bakhtiari, K. Daneshjou, " Low-Thrust Trajectory Design from LEO to GEO Using Reinforcement Learning", *Journal of Space Science, Technology and Applications*, vol 4 (1), p.: 64-77, 2024.

## Low-Trust Trajectory Design from LEO to GEO Using Reinforcement Learning

Hamed Soleymani, Majid Bakhtiari, Kamran Daneshjou

- \* PhD student in aerospace engineering at Iran University of Science and Technology, Tehran, Iran, [h\\_soleymani99@nt.iust.ac.ir](mailto:h_soleymani99@nt.iust.ac.ir) (Corresponding Author)
- School of Advanced Technologies - Iran University of Science and Technology, Tehran, Iran, [bakhtiari\\_m@iust.ac.ir](mailto:bakhtiari_m@iust.ac.ir)
- Faculty member of Iran University of Science and Technology, Tehran, Iran, [kjoo@iust.ac.ir](mailto:kjoo@iust.ac.ir)

### Abstract

In this paper, In the preliminary phase of space mission design, the selection of the spacecraft's trajectory is critical. This study simulates the dynamics of low-thrust orbital transfers within a two-dimensional orbital plane, employing a set of ordinary differential equations to represent the continuum of the spacecraft's orbital elements. These elements are encapsulated by six orbital parameters, manipulated under a defined thrust vector strategy within the action space, adhering to a specified policy framework. An agent, trained via a reinforcement learning algorithm within an actor-critic network architecture, is tasked with executing a low-thrust transfer between Low Earth Orbit (LEO) and Geostationary Orbit (GEO). The algorithm dynamically adjusts the spacecraft's trajectory, informed by initial orbital conditions and mission-specific constraints, to derive an optimal thrust angle trajectory and corresponding adjustments in orbital elements for the maneuver. To validate the algorithm's efficacy and robustness, a comparative analysis is conducted by implementing an alternative transfer mode at a varied orbital altitude. Additionally, the study explores the impact of adjusting the algorithm's degradation coefficient hyperparameter on the learning efficacy. Conclusively, the findings suggest that the agent, once adequately trained within the specified dynamical model, is capable of autonomously executing analogous orbital transfers. This is achieved without necessitating reiteration of the dynamical simulations, contingent solely upon the stipulation of initial and terminal orbital parameters



دسترس پذیر در نشانی:  
Journal.isrc.ac.ir

دو فصلنامه  
علوم، فناوری و  
کاربردهای فضایی

سال چهارم، شماره ۱، صفحه ۶۴-۷۷  
بهار و تابستان ۱۴۰۳

DOI:  
10.22034/jsssta.2024.446359.1153

تاریخچه داوری

دریافت: ۱۴۰۲/۱۲/۱۳

پذیرش: ۱۴۰۳/۰۲/۲۹

واژه‌های کلیدی

تراست پایین، عناصر مداری اعتدالی،  
یادگیری تقویتی، عامل، شبکه بازیگر،  
منتقد

نحوه استناد به این مقاله

حامد سلیمانی، مجید بختیاری، کامران  
دانشجو، "طراحی مسیر فضاپیمای  
تراست پایین از مدار لئو به ژئو با  
استفاده از یادگیری تقویتی"،  
دوفصلنامه علوم، فناوری و کاربردهای  
فضایی، جلد چهارم، شماره اول،  
صفحات ۶۴-۷۷، ۱۴۰۳.

مقاله پژوهشی

## طراحی مسیر فضاپیمای تراست پایین از مدار لئو به ژئو با استفاده از یادگیری تقویتی

حامد سلیمانی\*<sup>۱</sup>، مجید بختیاری<sup>۲</sup>، کامران دانشجو<sup>۳</sup>

<sup>۱</sup> - دانشجوی دکترا هوافضا دانشگاه علم و صنعت ایران، تهران، ایران h\_soleymani99@nt.iust.ac.ir

<sup>۲</sup> - عضو هیئت علمی دانشگاه علم و صنعت ایران، تهران، ایران bakhtiari\_m@iust.ac.ir

<sup>۳</sup> - عضو هیئت علمی دانشگاه علم و صنعت ایران، تهران، ایران kjoo@iust.ac.ir

\* نویسنده مسئول

### چکیده

در فاز اولیه طراحی مأموریت‌های فضایی، انتخاب دقیق مسیر فضاپیما از اهمیت بالاتری برخوردار است. در این پژوهش، دینامیک انتقال مداری تراست پایین دایروی صفحه‌ای بر اساس معادلات دیفرانسیل اعتدالی به‌عنوان محیط پیوسته برای متغیرهای مسئله که شش عنصر مداری اعتدالی یک فضاپیما هستند، شبیه‌سازی می‌شود. بردار رانش به‌عنوان فضای عمل تعریف شده و تحت یک سیاست انتخاب و به محیط اعمال می‌شود. عامل توسط الگوریتم یادگیری تقویتی شبکه بازیگر-منتقد برای انجام انتقال مداری تراست پایین از مدار لئو به مدار ژئو آموزش داده می‌شود. مسیر فضاپیما توسط الگوریتم مطابق با شرایط اولیه و قیود مأموریت جستجو می‌شود و در نهایت، پروفیل زاویه تراست مطلوب و تغییرات عناصر مداری مرتبط برای یک حالت مانور انتقال مداری به دست خواهند آمد. برای بررسی دقت و اعتبار الگوریتم در نتایج حالت اول مانور مداری، حالت دوم با اندازه تراست متفاوت پیاده‌سازی می‌شود. همچنین تأثیر تغییر فرایمتر ضریب تنزل الگوریتم بر روند یادگیری نیز بررسی می‌شود. در نهایت با در نظر گرفتن نتایج، عامل آموزش دیده در محیط دینامیک مسئله می‌تواند مأموریت‌های مشابه را بدون نیاز به شبیه‌سازی مجدد دینامیک مسئله و پارامترهای آن و فقط با تعیین شرایط اولیه و نهایی، با موفقیت به انجام برساند

## ۱- مقدمه

مسئله طراحی مسیرهای فضایی بهینه برای کاهش مصرف سوخت و زمان انتقال، موضوع اساسی در حوزه هوافضا است. به منظور دستیابی به مسیرهای بهینه سازگار با پارامترهای کلی مأموریت، مانند سیستم پیشران و جرم سوخت در دسترس، باید یک مسئله بهینه‌سازی جامع و پیچیده حل شود. این مسئله از دهه ۱۹۵۰ مورد توجه قرار گرفته و توسط پیشگامانی همچون لاودن، میل، کان وی، لیتن و ادلبام بررسی شده است. مسئله انتقال مداری تراست پایین به دلیل پیچیدگی خاصی که دارد، نیاز به روش‌های عددی دارد و به طور تحلیلی قابل حل نمی‌باشد. استفاده از پیشران سوزش پیوسته به طور گسترده‌ای برای ملاقات مداری<sup>[۱]</sup>، نگهداری مدار<sup>[۲]</sup>، انتقال مدار و مأموریت‌های فضایی بین سیاره‌ای مورد مطالعه قرار گرفته است. سیستم‌های پیشران تراست پایین می‌تواند امکان‌پذیری‌های مأموریت را با استفاده کارآمدتر از سوخت، بسیار بهبود بخشد که به نوبه خود اجازه افزایش نسبت بار محموله به وزن کل فضاپیما و یا کاهش قابل توجه حجم فضاپیما و در نتیجه ابعاد و سازه پرتابگر را می‌دهد. هنگام استفاده از پیشران تراست پایین برای انجام انتقالات مداری، ضربه ویژه<sup>[۳]</sup> بالا ترکیب شده با یک نیروی ویژه کوچک منجر به زمان انتقال بسیار طولانی می‌شود.<sup>[۱]</sup>

برای یافتن مسیر بهینه، نیاز به حل مسئله پیچیده ای است که «مسئله دو مقدار مرزی<sup>[۴]</sup>» نامیده می‌شود. روش‌های کلاسیک برای حل مسئله کنترل بهینه به طور کلی به سه دسته روش‌های مستقیم، غیر مستقیم و ترکیبی تقسیم می‌شوند. در روش‌های مستقیم، با گسسته کردن مسیر به قوس‌های کوتاه با قدر و جهت رانش ثابت و اعمال شرایط مرزی مناسب، مسئله پیوسته به یک مسئله برنامه‌ریزی غیرخطی تبدیل می‌شود.<sup>[۷]</sup> روش‌های غیرمستقیم از تکنیکی استفاده می‌کنند که به عنوان اصل حداقل پونتراگین<sup>[۵]</sup> شناخته می‌شود، که مسئله را به یک مسئله ارزش مرزی دو نقطه‌ای تبدیل می‌کند. از سوی دیگر روش‌های ترکیبی، بهره‌گیری از مزایای هر دو روش ذکر شده است. در سال‌های گذشته مطالعات زیادی در زمینه بهینه‌سازی زمان انتقال و مصرف سوخت (مینیمم کردن مقدار تغییرات سرعت) در انتقال‌های ضربه‌ای انجام شده است<sup>[۱۸ و ۱۹]</sup>. در

حوزه تراست پایین، انتقال حداقل زمان ژئوسنکرون برای فضاپیماهای تراست پایین با استفاده از روش‌های پیوسته توسط کایلو و همکاران<sup>[۲]</sup> بررسی شد. آن‌ها مسیرهای انتقال چندگانه چرخه سوخت را با استفاده از روش تک تیراندازی و روش هیبریدی پاول بهینه کردند. میچا کیم<sup>[۳]</sup> یک رویکرد جدید برای بهینه‌سازی مسیرهای تراست پایین پیوسته، با تمرکز بر یک الگوریتم بهینه‌سازی مبتنی بر روش‌های غیرمستقیم پیشنهاد کرد و یک الگوریتم تبرید تطبیقی<sup>[۴]</sup> را برای تولید حدس اولیه معرفی کرد. رایان راسل<sup>[۴]</sup> از حساب تغییرات و نظریه بردار آغازگر برای انتقال زمین-ماه با رانش کم با چرخش‌های متعدد استفاده کرد. در این کار، با استفاده از جستجوی جهانی برای راه‌حل‌های روش غیرمستقیم محلی، به یک جبهه بهینه پرتو برای زمان پرواز و جرم نهایی دست یافت. در سال ۲۰۱۳ یک مطالعه مقایسه‌ای بر اساس روش‌های غیرمستقیم برای حل مسئله کنترل بهینه غیرخطی برای مسیر تراست پایین فضایی توسط نوبی و همکاران<sup>[۱۶]</sup> ارائه شد. که لی هوانگ<sup>[۵]</sup> بر بهینه‌سازی انتقال‌های تراست پایین به نقاط لاگرانژی تمرکز کرد و از برنامه‌ریزی درجه دوم متوالی برای بهینه‌سازی مسیرهای انتقال از مدارهای پارکینگ به مدارهای لاگرانژی استفاده کرد. بروس کانوی و مارو پونتانی<sup>[۶]</sup> یک روش بهینه‌سازی ازدحام ذرات<sup>[۶]</sup> غیرمستقیم را پیشنهاد کردند، که پیش‌نیازهای ضروری بهینه‌سازی را با یک الگوریتم اکتشافی یکپارچه می‌کند. جاناتان عزیز و همکاران<sup>[۸]</sup> یک تبدیل سلندمن<sup>[۶]</sup> با برنامه‌ریزی پویای دیفرانسیلی<sup>[۶]</sup> را برای بهینه‌سازی معادلات حرکت فضاپیما اعمال کرد و متغیر مستقل را به زاویه مدار برای بهینه‌سازی در اطراف اجسام سیاره‌ای تغییر داد. ویلیامز و کاورستون-کارول<sup>[۹]</sup> با استفاده از الگوریتم ژنتیک<sup>[۶]</sup>، مسئله حداقل سوخت را برای فضاپیماهای پیشران الکتریکی خورشیدی مطالعه کردند. فکور و همکاران<sup>[۲۰]</sup> در سال ۲۰۱۹ مسئله کنترل بهینه مسیر تراست پایین دایروی از مدار لثو به ژئو را با ترکیب روش تحلیلی و الگوریتم جمعیتی کلونی زنبور مصنوعی<sup>[۶]</sup> برای حالت‌های کمترین زمان، کمترین تلاش کنترلی مطالعه کرده‌اند. در<sup>[۱۰]</sup>، ژنوب وانگ یک استراتژی بهینه‌سازی محدب برای حل مسائل انتقال بهینه تراست پایین، با تاکید بر بهره‌وری سوخت، توسعه داد. در واقع هر دو روش مستقیم و غیر مستقیم

مداری به ارتفاع بالا، از الگوریتم شبکه بازیگر-منتقد<sup>[۱۱]</sup> استفاده می شود. دینامیک انتقال مداری تراست پایین دایروی هم صفحه بر اساس معادلات دیفرانسیل اعتدالی به عنوان محیطی برای تعامل عامل که فضایی پیوسته برای متغیرهای مسئله که شش عنصر مداری اعتدالی یک فضاپیما هستند، شبیه سازی می شود. برادر رانش به عنوان فضای عمل تعریف شده و تحت یک سیاست تصادفی انتخاب و به محیط اعمال می شود. عامل توسط الگوریتم شبکه بازیگر-منتقد آموزش داده می شود تا بتواند انتقال مداری تراست پایین مورد نظر را انجام دهد. مسیر بهینه توسط این الگوریتم مطابق با شرایط اولیه و قیود ماموریت (مانند اندازه نیروی تراست) جستجو می شود و در نهایت پروفیل های زاویه تراست بهینه برای دو حالت مانور انتقال مداری و یک حالت مانور تغییر زاویه میل مدار به دست خواهند آمد.

## ۲- یادگیری ماشین؛ یادگیری تقویتی

هر مقاله یادگیری ماشین به طور کلی به سه دسته تقسیم می شود: یادگیری نظارت شده<sup>[۱۲]</sup>، یادگیری بدون نظارت<sup>[۱۳]</sup> و یادگیری تقویتی. آخرین دسته، یادگیری تقویتی، که هدف این پژوهش است، به یادگیری ماشین اجازه می دهد تا برای مسائلی که هیچ داده ای در دسترس نیست استفاده شود. یادگیری تقویتی عیوب کلی روش های قبلی را برطرف کرده و بر روی طیف وسیعی از مسائل کنترلی پیاده سازی شده است. سیستم یادگیری در یادگیری ماشینی عامل<sup>[۱۴]</sup> نامیده می شود و رفتار بهینه آن عامل با یک نگاهت از ناحیه حالت  $X$  به ناحیه عمل  $A$  به صورت  $A: X \rightarrow S$  تعریف می شود. سیاست<sup>[۱۵]</sup> یک اصطلاح کلی است که یک عامل برای انجام یک عمل در یک محیط از آن پیروی می کند. مقدار پاداش<sup>[۱۶]</sup> عددی است که به یک عامل می گوید که یک عمل در یک انتقال حالت چقدر خوب است. [۱۱]

فرآیند تصمیم گیری مارکوف چارچوبی است که برای تعریف مسائل یادگیری تقویتی استفاده می شود. در یادگیری تقویتی، یک عامل در محیطی با هدف به حداکثر رساندن ارزش پاداش قرار می گیرد. عامل با انجام یک عمل تعیین شده توسط یک سیاست با محیط تعامل دارد. وقتی حالت خاصی حاصل شد، پاداشی توسط محیط به عامل داده می شود. ارزش<sup>[۱۷]</sup> به عنوان بازگشت مورد انتظار<sup>[۱۸]</sup> در حالت  $S$  در هنگام پیروی از یک خط

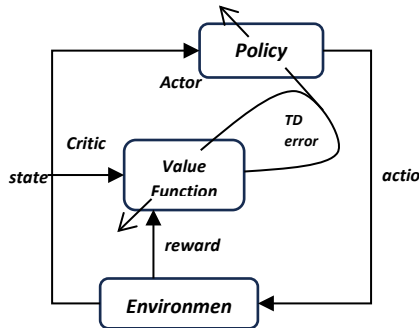
ماهیت محلی دارند. این بدان معنی است که راه حل بهینه ای که آن ها می توانند شناسایی کنند به نخستین راه حل موجود بستگی دارد. نتیجه عددی به دست آمده از تکنیک های قطعی به طور معمول در مجاورت حدس است [۶]. در نتیجه بررسی روش های حل مسئله کنترل بهینه که در بالا بیان شد، می توان گفت که انتخاب بین این روش ها به نیازهای خاص مسئله از جمله پیچیدگی قیود، الزامات دقت عددی و امکان سنجی حدس اولیه بستگی دارد. همچنین روش های مبتنی بر جمعیت مانند الگوریتم ازدحام ذرات و ژنتیک مسائل مقادیر اولیه<sup>[۱۹]</sup> هستند و برای همگرایی نیاز به راه حل اولیه خوبی دارند. از آنجایی که مسئله انتقال مدار تراست پایین یک مسئله مقدار مرزی دونقطه ای است و رسیدن به حالت نهایی مطلوب مهم است، استفاده از الگوریتمی لازم است که بتواند این فرآیند را برآورده کند. علاوه بر این، الگوریتم ازدحام ذرات و ژنتیک می توانند در حداقل های محلی به خصوص در مسائل دارای فضاهای حالت بسیار پیچیده چند حالتی به دام بی افتند، عملکرد آن ها می تولد با افزایش ابعاد پیچیدگی مسئله تضعیف شود، به تنظیمات پارامترهایشان حساس هستند.

این اشکالات تقریباً توسط روش های یادگیری ماشینی پوشش داده می شوند و باعث می شوند که این روش ها در مسائل ابعاد بالا و بسیار پیچیده استفاده شوند، در حداقل های محلی به دام نیافتند و در کاربردهای خودکار و خودران استفاده شوند. یکی از شاخه های مهم در یادگیری ماشینی، یادگیری تقویتی<sup>[۲۰]</sup> است که در سیستم های مختلف مانند ماشین ها و ربات های صنعتی استفاده می شود. این نوع یادگیری به دلیل ماهیت تعاملی اش، بیشتر شبیه به یادگیری از محیط اطراف است که در انسان و حیوانات دیده می شود. در زمینه یادگیری تقویتی، عامل برای دستیابی به اهداف خود نیازی به دانش قبلی از دینامیک محیط ندارد. این ویژگی، الگوریتم های یادگیری تقویتی را قادر می سازد تا به طور مؤثری با مسائل مدل سازی تحلیلی پیچیده مقابله کنند. از این رو، این ویژگی امکان آموزش کارآمد عامل را فراهم می کند و آن را برای استفاده در موقعیت های مختلف با دینامیک مسئله مشابه قابل تطبیق می کند.

در این پژوهش برای تحلیل عملکرد این نوع الگوریتم های یادگیری ماشینی به عنوان یک کار پیشگام در زمینه انتقالات

یک تابع ارزش ارزیابی می‌شود. الگوریتم‌گرادین سیاست قطعی عمیق  $Q^*$  یکی از چندین الگوریتم‌های بازیگر-منتقد است. به طور معمول، منتقدیک تابع ارزش حالت  $Q^*$  است. با انتخاب شدن یک عمل، منتقد وضعیت جدید را ارزیابی می‌کند تا مشخص کند که آیا اوضاع بهتر یا بدتر از حد انتظار پیش رفته است. این ارزیابی خطای اختلاف زمانی است که به آن خطای تفاضل زمانی  $TD$  می‌گویند: [۱۱]

$$\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t) \quad (2)$$



شکل ۲. معماری شبکه بازیگر-منتقد [۱۲]

که در آن  $V(s_t)$  تابع ارزش فعلی است که توسط منتقد پیاده‌سازی شده است. این خطای تفاضل زمانی می‌تواند برای ارزیابی عمل انتخاب شده  $a_t$ ، عمل انجام شده در حالت  $s_t$  استفاده شود. اگر خطای تفاضل زمانی مثبت باشد، نشان می‌دهد که تمایل به انتخاب  $a_t$  باید برای آینده تقویت شود، در حالی که اگر خطای تفاضل زمانی منفی باشد، نشان می‌دهد که این تمایل باید ضعیف شود. [۱۱]

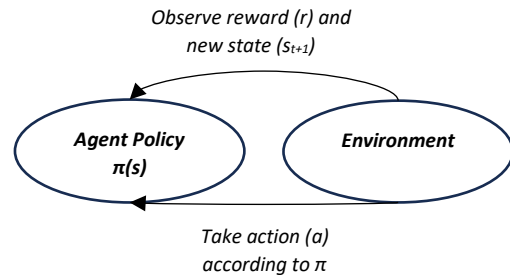
مدل منتقد تحت عنوان یک شبکه عصبی نمایش داده می‌شود که ورودی‌های آن حالت و عمل یک عامل و خروجی آن یک مقدار  $q$  منفرد است. [۱۱]

یک شبکه عصبی پیش‌خور از لایه‌هایی تشکیل شده است که از طریق وزن‌ها به یکدیگر متصل شده‌اند. به‌طور کلی از یک لایه ورودی، یک یا چند لایه پنهان و یک لایه خروجی تشکیل شده است که در شکل ۳ نشان داده شده است.

مشی  $\pi$  تعریف می‌شود. تابع ارزش استخراج شده و مشتق شده از معادله بلمن، کل پاداش دریافتی با انجام یک عمل در یک حالت خاص را بیان می‌کند. از نظر یک فرآیند مارکوف محدود، تابع ارزش را می‌توان به صورت زیر تعریف کرد:

$$v_\pi = E_\pi[G_t | S_t = s] = E_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right] \quad (1)$$

فرآیند تصمیم‌گیری مارکوف در شکل (۱) نمایش داده شده است. [۱۱]



شکل ۱. فرآیند تصمیم‌گیری مارکوف [۱۱]

$E_\pi$  مقدار مورد انتظار یک متغیر تصادفی است زمانی که یک عامل از یک سیاست خاص  $\pi$  پیروی می‌کند؛  $R$  مقدار پاداش در یک حالت در زمان  $t$  است و  $\gamma$  یک ضریب تنزل  $Q^*$  است که مقدار آن بین صفر و ۱ است. ضریب تنزل نشانه از میزان تأثیر پاداش‌های آینده بر تابع ارزش است. [۱۱]

گرادین سیاست  $Q^*$  یک عمل را بر اساس یک سیاست پارامتری شده بدون یادگیری تابع ارزش انتخاب می‌کند. با روش‌های گرادین سیاست، سیاست باید پیوسته و قابل انتگرال‌گیری باشد. این اجازه می‌دهد تا روش‌های گرادین سیاست در فضاهای عمل پیوسته اعمال شوند. دو الگوریتم اصلی که گرادین‌های سیاست را اعمال می‌کنند عبارتند از REINFORCE، یک گرادین سیاست مبتنی بر مونت کارلو و روش‌های Actor-Critic. این الگوریتم‌ها را می‌توان در فضاهای عمل پیوسته اعمال کرد. [۱۱]

الگوریتم‌های منتقد بازیگر نشان داده شده در شکل ۲، روش‌های تفاضل زمانی و گرادین سیاست را برای کاهش اشکالات دو نوع الگوریتم ترکیب می‌کنند. الگوریتم‌های بازیگر-منتقد شامل بازیگری است که عملی را بر اساس یک سیاست انتخاب می‌کند در حالی که سیاست توسط منتقد با استفاده از

$$k = \tan\left(\frac{i}{2}\right) \sin(\Omega) \quad (8)$$

$$L = \Omega + \omega + \theta \quad (9)$$

حال معادلات دیفرانسیل معمولی مرتبه اول گاوس برحسب

[۱۱] عناصر اعتدالی بازنویسی می‌شود:

$$\frac{dp}{dt} = \frac{2p}{w} \sqrt{\frac{p}{\mu}} F_s \quad (10)$$

$$\frac{df}{dt} = \sqrt{\frac{p}{\mu}} \left[ \sin L F_r + [(w+1) \cos L + f] \frac{F_s}{w} - (h \sin L - k \cos L) \frac{g F_w}{w} \right] \quad (11)$$

$$\frac{dg}{dt} = \sqrt{\frac{p}{\mu}} \left[ -\cos L F_r + [(w+1) \sin L + g] \frac{F_s}{w} + (h \sin L - k \cos L) \frac{g F_w}{w} \right] \quad (12)$$

$$\frac{dh}{dt} = \sqrt{\frac{p}{\mu}} \frac{s^2 F_w}{2w} \cos L \quad (13)$$

$$\frac{dk}{dt} = \sqrt{\frac{p}{\mu}} \frac{s^2 F_w}{2w} \cos L \quad (14)$$

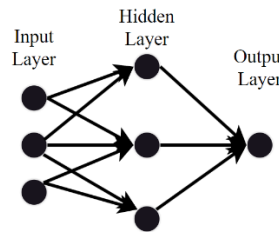
$$\frac{dL}{dt} = \sqrt{\mu p} \left(\frac{w}{p}\right)^2 + \frac{1}{w} \sqrt{\frac{p}{\mu}} (h \sin L - k \cos L) F_w \quad (15)$$

که در آن  $F_s$ ،  $F_r$  و  $F_w$  به ترتیب نیرو تراست در جهت‌های شعاعی، مماسی و نرمال بر صفحه مداری هستند. جایی که  $r$  در امتداد جهت شعاعی مدار،  $s$  جهت محیطی و  $w$  در امتداد بردار ممتموم زاویه‌ای است. [۱۱]

$$s^2 = 1 + h^2 + k^2 \quad (16)$$

$$w = 1 + f \cos L + g \sin L \quad (17)$$

در این بررسی برای تحلیل عملکرد این نوع از الگوریتم‌های یادگیری ماشینی و ایجاد یک پژوهش اولیه و هم‌چنین به‌عنوان یک کار اولیه در زمینه انتقال مداری ارتفاع‌های بالا با این دسته از روش‌ها برای محققان دیگر، این پیاده‌سازی از حالت بدون مدل استفاده می‌شود. در اینجا موارد مهم و مدنظر برای انجام فرآیند طراحی مسیر تغییرات نیم قطر اصلی، خروج از مرکز و تغییرات زاویه نیروی رانش رانشگر است.



شکل ۳. شبکه عصبی پیش‌خور [۱۱]

### ۳- دینامیک مسئله انتقال مداری تراست پایین

کنترل یک فضاپیما را می‌توان با فرآیند تصمیم‌گیری مارکوف توصیف کرد، به این معنا که در هر مرحله زمانی فضاپیما حالت خاصی دارد و باید تصمیم گرفت که چه اقدامی انجام دهد تا در نهایت به یک هدف برسد. [۱۳]  
در مانورهای تراست پایین، یک نیروی خارجی  $F$  به فضاپیما معادله حرکت اضافه می‌شود: [۱۴]

$$\ddot{\mathbf{r}} = \frac{\mu}{|\mathbf{r}|^3} \mathbf{r} + \frac{\mathbf{F}}{m} \quad (3)$$

جایی که  $m$  جرم فضاپیما و  $\mu$  ثابت گرانشی استاندارد جسم مرکزی است. [۱۴]

دینامیک غیرخطی مانورهای تراست پایین را می‌توان با استفاده از معادلات تغییرات گاوس [۱۵] توسط معادلات برحسب عناصر مداری کپلری بیان کرد که در مرجع ۶ بیان شده است. اما یکی از معایب استفاده از عناصر مداری کپلری، تکینگی در شب مداری یا خروج از مرکز صفر است. لذا به سراغ عناصر مداری اعتدالی [۱۶] می‌رویم. مجموعه‌ای از عناصر مداری هستند که برای تحلیل و طراحی مسیر بهینه مفید هستند. آن‌ها برای مدارهای دایروی، بیضی و هذلولی معتبر هستند. این معادلات اعتدالی اصلاح‌شده، هیچ تکینگی برای خروج از مرکز صفر و زاویه شیب مداری برابر با ۰ و ۹۰ درجه را نشان نمی‌دهند و می‌توان از این تکینگی‌ها جلوگیری کرد.

$$p = a(1 - e^2) \quad (4)$$

$$f = e \cos(\omega + \theta) \quad (5)$$

$$g = e \sin(\omega + \theta) \quad (6)$$

$$h = \tan\left(\frac{i}{2}\right) \cos(\Omega) \quad (7)$$

#### ۴- نتایج

برای فضاپیما، تابع ارزش تابعی از تفاضل بین حالت فعلی و نهایی و هدف فضاپیما است. در محیط یادگیری تقویتی یک الگوریتم بر مبنای مدل بازیگر-منتقد، عناصر مداری اعتدالی ۶ گانه به عنوان پارامترهای دخیل در الگوریتم یا به عبارت بهتر حالت‌های مسئله خواهند بود. هدف این است که این عامل با تعامل با محیط بتواند بهترین دنباله از همین حالت‌ها را پیدا کند که یک مسیر برای انتقال مداری از مدارهای پایین به مدارهای بالاتر به حساب آید. در بخش شبکه بازیگر کار به این صورت است که بر اساس بازه‌ای که برای مقدار بزرگی نیروی رانش که توسط کاربر بر مبنای قیود و الزامات مأموریت انتخاب و یک دسته عمل به عامل داده می‌شود. عامل به صورت تصادفی از بین آن‌ها عملی را انتخاب کرده و به محیط مسئله اعمال می‌کند. از سوی دیگر شبکه منتقد که وظیفه ارزیابی تابع ارزش حالت مسئله را دارد، در هر اپیزود فرآیند آموزش، با ارزیابی‌هایی که از مقدار ارزش هر حالت انجام شده بر اساس عمل انتخاب شده توسط سیاست عامل، سعی می‌کند مقادیر به دست آمده به مقدار پاداش میانگین نزدیک‌تر کند و در نهایت نیز با برقراری شرایط توقف و رسیدن به نقطه نهایی با توجه به شرایط مرزی تعیین شده، الگوریتم را متوقف می‌کند و در نهایت پروفیل زاویه تراست بهینه برای انجام این انتقال حاصل خواهد شد. در این بین شرط توقف مسئله با توجه به اینکه مأموریت مانور افزایش مدار باشد یا مانوری برای تغییر خروج از مرکز یا زاویه میل تعیین می‌شود. شکل ۶ پیاده سازی الگوریتم را در نرم افزار نمایش می‌دهد.

لازم به ذکر است که با توجه به ساده و ملموس بودن عناصر مداری کپلری در مباحث مکانیک مداری، و همچنین بیان بهتر نتایج خروجی حاصل از الگوریتم در انتقال مداری مورد نظر، خروجی‌ها هم در قالب عناصر مداری اعتدالی و هم کپلری نمایش داده خواهد شد.

تابع پاداش (مثبت برای یک پاداش، منفی برای یک جریمه) پاداش‌هایی را برای قرار گرفتن در یک حالت خاص یا برای انجام عملی در یک حالت خاص مشخص می‌کند. تابع پاداش به طور ضمنی هدف یادگیری را بیان می‌کند. به گونه‌ای که تابع پاداش به اینکه چطور باید سیستم (یعنی فرآیند تصمیم‌گیری مارکوف) کنترل شود سمت‌وسو می‌دهد. [۱۵]

تابع پاداش می‌تواند یکی از چالش‌برانگیزترین انتخاب‌ها در هنگام طراحی یک محیط باشد. کیفیت تابع پاداش می‌تواند تعیین کند که عامل چقدر سریع به یک راه حل همگرا می‌شود. دادن یک پاداش پراکنده به عامل ممکن است زمان اجرا را به میزان قابل توجهی افزایش دهد و همگرایی عامل را به یک راه حل دشوار کند.

تابع پاداش در این تحقیق بر اساس اختلاف مجذور بین حالت هدف و حالت فعلی است که در زیر نشان داده شده است:

$$r_p = \frac{\sqrt{(p_{targ} - p(t))^2}}{p_{targ}} \quad (18)$$

$$r_f = \frac{\sqrt{(f_{targ} - f(t))^2}}{f_{targ}} \quad (19)$$

$$r_g = \frac{\sqrt{(g_{targ} - g(t))^2}}{g_{targ}} \quad (20)$$

$$r_h = \frac{\sqrt{(h_{targ} - h(t))^2}}{h_{targ}} \quad (21)$$

$$r_k = \frac{\sqrt{(k_{targ} - k(t))^2}}{k_{targ}} \quad (22)$$

$$r = -(\alpha_p r_p + \alpha_f r_f + \alpha_g r_g + \alpha_h r_h + \alpha_k r_k) \quad (23)$$

ضرایب  $\alpha_p, \alpha_f, \alpha_g, \alpha_h, \alpha_k$  برای مقیاس بندی وزن نسبی عناصر مداری هستند. برای اینکه عامل یک سیاست خوب را یاد بگیرد، مقیاس بندی مناسب تابع پاداش لازم است. با توجه به تعریف تابع پاداش و اینکه مانور مورد نظر افزایش ارتفاع یا اصلاح زاویه میل مدار باشد، ضرایب مربوط به عناصر مداری اعتدالی مؤثر در مأموریت باید بالاتر از سایر ضرایب تنظیم شوند، زیرا تابع پاداش به طور ضمنی هدف یادگیری را بیان می‌کند. از آنجایی که عموماً در مانورهای انتقال مدار هدف رسیدن به ارتفاع و جهت گیری مشخص شده مدار نهایی است و فقط حالت مداری در نظر گرفته می‌شود، عنصر مداری آنومالی حقیقی را می‌توان حذف کرد. بنابراین، عنصر مداری زاویه آنومالی حقیقی  $\theta$  در هیچ‌یک از مأموریت‌ها هدف قرار نمی‌گیرد یا در تابع پاداش گنجانده نمی‌شود. این ضرایب در جدول ۱ ارائه شده است.

جدول ۱. ضرایب تابع پاداش و مقادیر نوسان مجاز عناصر مداری

مقدار نوسان مجاز عناصر مداری		ضرایب تابع پاداش	
۱۰	p(km)	۱۰	$\alpha_p$
۰.۰۱	f(rad)	۵	$\alpha_f$
۰.۰۱	g(rad)	۵	$\alpha_g$
۰.۰۰۱	h(rad)	۱	$\alpha_h$
۰.۰۰۱	k(rad)	۱	$\alpha_k$

جدول ۲. مشخصات حالت اول مانور افزایش مدار

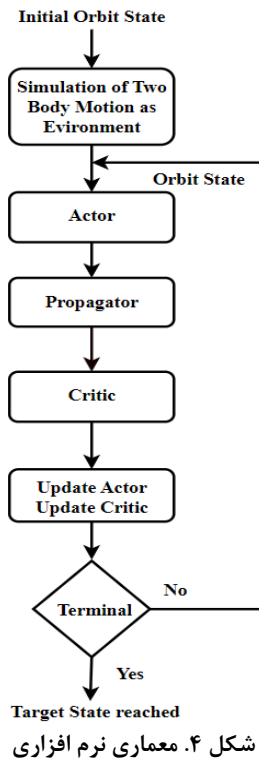
حالت دوم مانور مداری	حالت اول مانور مداری	
۳۰۰۰	۳۰۰۰	جرم فضاییما(kg)
۴	۳	مقدار تراست(N)
۱۰۰۰	۱۰۰۰	ارتفاع مدار اولیه(km)
۳۵۷۸۶	۳۵۷۸۶	ارتفاع مدار نهایی(km)

با در نظر گرفتن به مقدار نوسان مجاز عناصر مداری، جدول ۴ به وضوح نتایج به دست آمده از الگوریتم در انتقال تراست پایین را نشان می دهد که به طور منطقی قابل قبول هستند. برای هر دو حالت مانور انتقال مداری، پارامترهای شبکه های عصبی مطابق جدول ۵ تنظیم شده است.

جدول ۳. مقادیر عناصر مداری حاصل از الگوریتم یادگیری

تقویتی در حالت اول مانور افزایش مدار

زمان (انتقال) (روز)	عناصر مداری اعتدالی			عناصر مداری کپلری		
	اولیه	نهایی		اولیه	نهایی	
۴۹.۴۱	7.378 1e+0 3	4.216 8e+0 4	p = a	7.378 1e+0 3	4.216 8e+0 4	a(km)
	0	-0.004 1	f(rad)	0.0	0.012 8	e
	0	0.011 2	g(rad)	0.0	0.0	i(deg)
	0	-4.992 8e-09	h(rad)	0.0	2.057 7e+0 2	$\omega$ (deg)
	0	2.174 5e-10	k(rad)	0.0	0.0	$\Omega$ (deg)



هدف اصلی که مانور افزایش مدار به مدار هم ارتفاع مدار ژئو است، در دو حالت با ارتفاع یکسان اما مقدار تراست متفاوت برای مدار اولیه به منظور بررسی عملکرد و دقت الگوریتم انجام می شود:

- ۱) مدار اولیه در ارتفاع ۱۰۰۰ کیلومتری از سطح زمین و تراست ۳ نیوتن
- ۲) مدار اولیه در ارتفاع ۱۰۰۰ کیلومتری از سطح زمین و تراست ۴ نیوتن

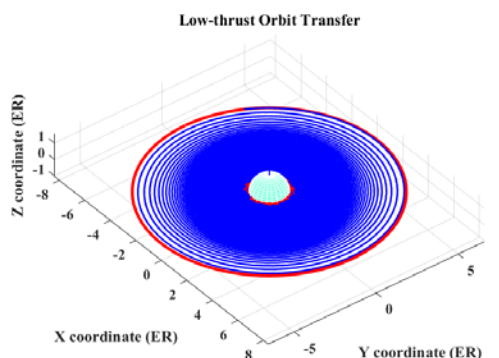
#### ۱) حالت نخست مانور افزایش مدار

مشخصات مانور افزایش مدار در جدول ۳ ارائه شده است. سپس با ایجاد تغییر در فرآیندهای الگوریتم، میزان تأثیر تغییرات آن ها بر فرآیند جستجو، بهره برداری و اهمیت دهی به پاداش های آینده بر مبنای ضرایب های تنزل، نرخ یادگیری و بررسی خواهد شد. لازم است مقادیر نوسان برای عناصر مداری مورد نظر ایجاد شود. هنگامی که تفاوت بین حالت مدار هدف و وضعیت مدار فعلی مطابق مقادیر نوسان مجاز مشخص شده در جدول ۲ باشد، اپیزود با موفقیت خاتمه می یابد.

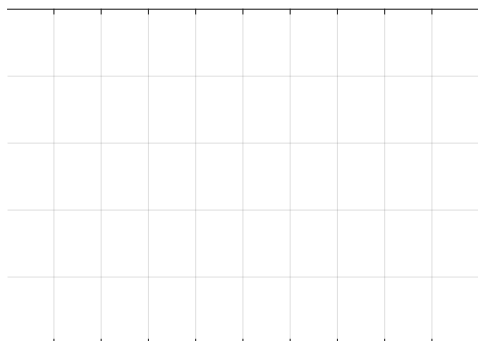
جدول ۴. پارامترهای شبکه عصبی بازیگر منتقد در مانور افزایش مدار

	شبکه بازیگر			شبکه منتقد		
	حالت اول	حالت دوم	حالت سوم	حالت اول	حالت دوم	حالت سوم
تعداد گره لایه پنهان اول	۲۸۰	۲۸۰	۲۸۰	۳۶۰	۳۶۰	۳۶۰
تعداد گره لایه پنهان دوم	۲۰۰	۲۰۰	۲۰۰	۲۸۰	۲۸۰	۲۸۰
تابع فعال ساز	ReLU	ReLU	ReLU	ReLU	ReLU	ReLU
نرخ یادگیری	۰.۰۰۱	۰.۰۰۱	۰.۰۰۱	۰.۰۱	۰.۰۱	۰.۰۱
ضریب تنزل	۰.۹	۰.۹۵	۰.۹۵	۰.۹۰	۰.۹۵	۰.۹۵

شکل ۶. تابع پاداش در حالت اول مانور افزایش مدار



شکل ۷. مدارات انتقال در حالت اول مانور افزایش مدار



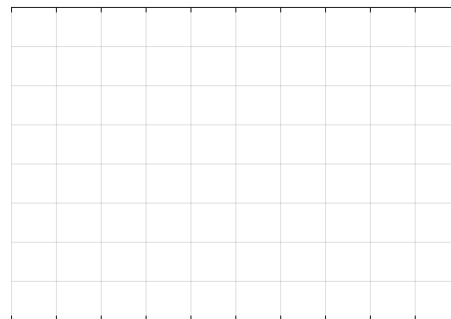
شکل ۸. پروفیل زاویه انحراف تراستر در حالت اول مانور افزایش مدار

در شکل ۶، نمودار زرد، تعاملات عامل در محیط و پاداش‌های جمع آوری شده از عامل را در اپیزودهای متوالی نشان می‌دهد. نوسان اولیه به این دلیل است که عامل سعی می‌کند بین فضای محیط را بیشتر مورد جستجو قرار دهد تا دنباله‌ای مطلوب برای حالت‌ها در این مأموریت پیدا کند. در نهایت، فرآیند با برآوردن شرایط توقف و رسیدن به حالت‌ها به حالت‌های مداری نهایی مطلوب به پایان می‌رسد. شکل ۷، مدارات انتقال در انجام این مأموریت را نشان می‌دهد که مدار قرمز رنگ همان مدار ژئو است که با موفقیت فضاییما به آن رسیده است. با توجه به فرضیات انتقال مداری در نظر گرفته شده، بردار تراست وارده در این انتقال مماس بر مسیر حرکت فضاییما به دست آمده است که در شکل ۸ نشان داده شده است. تغییرات نیم قطر اصلی بیضی و خروج از مرکز مدارات انتقال نیز به ترتیب در شکل‌های ۹ و ۱۰ نمایش داده شده است که می‌توان دید به ارتفاع مداری و خروج از مرکز مدار ژئو دست پیدا کرده است.

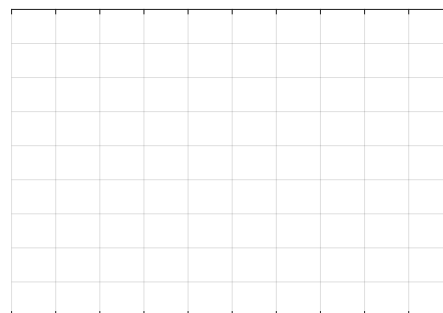
بر مسیر حرکت ماهواره اعمال شده است تا در نهایت شرایط نیم قطر اصلی و خروج از مرکز مدار ژئو مطابق شکل های ۱۴ و ۱۵ برآورده شود. مقدار خروج از مرکز در مدار نهایی همان طور که در شکل ۱۴ مشخص است به مقدار  $0.1176$  می رسد و همچنین نوسان آن در نهایت کنترل می شود که با توجه به تابع پاداش و هدف تعریف شده برای مسئله قابل قبول می باشد.

جدول ۵. مقادیر عناصر مداری حاصل از الگوریتم یادگیری تقویتی در حالت دوم مانور افزایش مدار

زمان انتقال (روز)	عناصر مداری اعتدالی			عناصر مداری کپلری		
	نهایی	اولیه		نهایی	اولیه	
۳۷.۱۰	$4.2168e+04$	$7.3781e+03$	$p = a$	$4.2167e+04$	$7.3781e+03$	$a(km)$
	$0.0041$	$0.0$	$f(rad)$	$0.0117$	$0.0$	$e$
	$0.0112$	$0.0$	$g(rad)$	$0.0$	$0.0$	$i(deg)$
	$4.9928e-09$	$0.0$	$h(rad)$	$1.7074e+02$	$0.0$	$\omega(deg)$
	$2.1745e-10$	$0.0$	$k(rad)$	$1.7844e+02$	$0.0$	$\Omega(deg)$



شکل ۹. تغییرات نیم قطر اصلی در حالت اول مانور مداری



شکل ۱۰. تغییرات خروج از مرکز مدارات در حالت اول مانور مداری

## ۲) حالت دوم مانور افزایش مدار

در حالت دوم، همانطور که در جدول ۵ مشخص است، مقدار ضریب تنزل به  $0.95$  افزایش داده می شود تا در ضمن انجام مأموریت، تأثیر پاداش های آینده بر عملکرد الگوریتم نیز بررسی شود. جدول ۶ به وضوح نتایج به دست آمده از الگوریتم در انتقال تراست پایین مربوطه را نشان می دهد که به طور منطقی قابل قبول هستند. نمودار پاداش حالت دوم در شکل ۱۱ قابل مشاهده است. با افزایش این ضریب، تمرکز و سرعت عامل بر پاداش های آینده بیشتر شده است و با توجه به اینکه در حالت اول یکبار عامل آموزش دیده است، فرآیند همگرایی و جمع آوری پاداش مثبت که به معنی مطلوب بودن حالت های به وجود آمده برای این انتقال مداری است، با سرعت و دقت بهتری پیش رفته است. در نهایت عامل با جمع آوری حداکثر پاداش مثبت با موفقیت به شرایط مدار ژئو دست پیدا می کند. در شکل ۱۲ مشاهده می شود که فضاپیما با موفقیت در مدار مطلوب قرار گرفته است. منحنی زاویه تراست نیز نشان می دهد که در تمام مسیر مماس

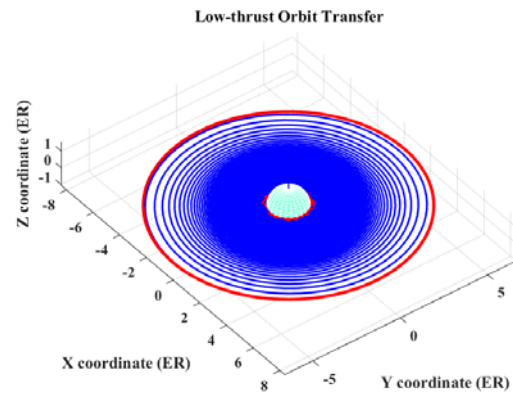
شکل ۱۱. نمودار پاداش در حالت دوم مانور افزایش ارتفاع مدار

### ۳) حالت سوم مانور افزایش مدار

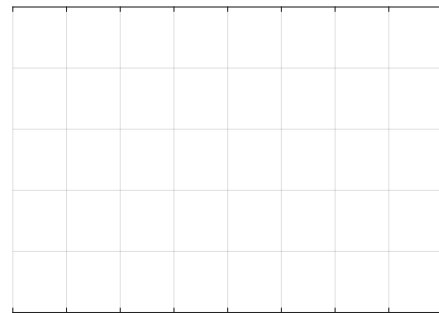
در حالت سوم با هدف بررسی دقت و صحت الگوریتم در نتایج، ارتفاع مداری به ۱۵۰۰ کیلومتری از سطح زمین افزایش یافته است. مشخصات شبکه ها و فرآپارمترهای الگوریتم طبق جدول ۵ مانند حالت دوم تنظیم شده است. مقادیر عناصر مداری مدار نهایی در جدول ۶ نمایش داده شده است. با توجه به اینکه عامل در دو حالت قبل با دینامیک مسئله یادگیری را انجام داده است، در این جا همانطور که از شکل ۱۶ مشخص است عامل با استفاده از تعاملات قبلی، با دریافت پاداش های مثبت و بدون داشتن تعاملات و پاداش منفی (شیب مثبت) توانسته است فرآیند یادگیری را با موفقیت و دقت بهتری نسبت به تو دو حالت قبل به انجام رساند و دنباله ای از حالت ها را برای مسیر تراست پایین مورد نظر مطابق شکل ۱۷ ترسیم نماید. مقادیر تغییرات زاویه تراست، نیم قطر اصلی و خروج از مرکز مدار های انتقال تراست پایین در شکل های ۱۸ تا ۲۰ نمایش داده شده است.

جدول ۶. مقادیر عناصر مداری حاصل از الگوریتم یادگیری تقویتی در حالت سوم مانور افزایش مدار

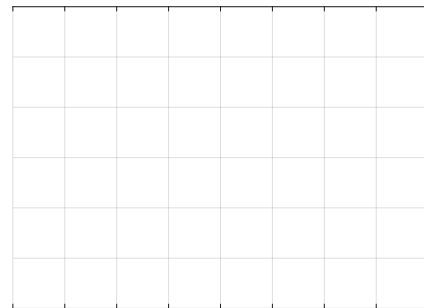
زمان انتقال (روز)	عناصر مداری اعتدالی			عناصر مداری کپلری		
	نهایی	اولیه		نهایی	اولیه	
۲۵.۰۵	4.216 8e+0 4	7.878 1e+0 3	p = a	4.216 7e+0 4	7.878 1e+0 3	a(km)
	- 0.004 1	0.0	f(rad)	0.011 7	0.0	e
	0.011 2	0.0	g(rad)	0.0	0.0	i(deg)
	- 4.992 8e-09	0.0	h(rad)	1.707 4e+0 2	0.0	$\omega$ (deg)
	2.174 5e-10	0.0	k(rad)	1.784 4e+0 2	0.0	$\Omega$ (deg)



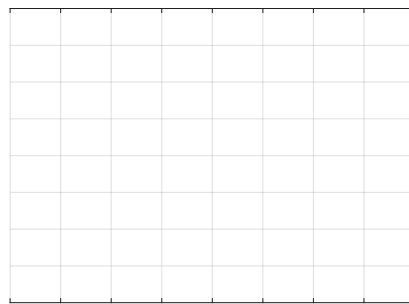
شکل ۱۲. مدارات انتقال در حالت دوم مانور افزایش مدار



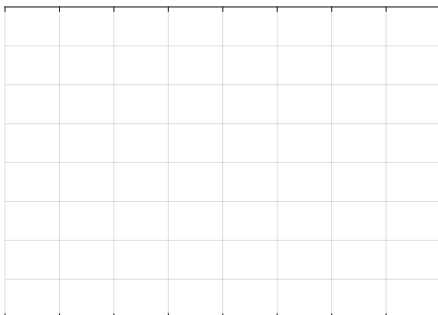
شکل ۱۳. پروفیل زاویه انحراف تراستر در حالت دوم مانور افزایش مدار



شکل ۱۴. تغییرات خروج از مرکز در حالت دوم مانور افزایش مدار



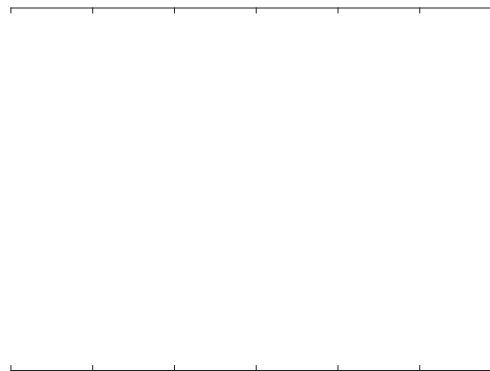
شکل ۱۵. تغییرات نیم قطر اصلی در حالت دوم مانور افزایش مدار



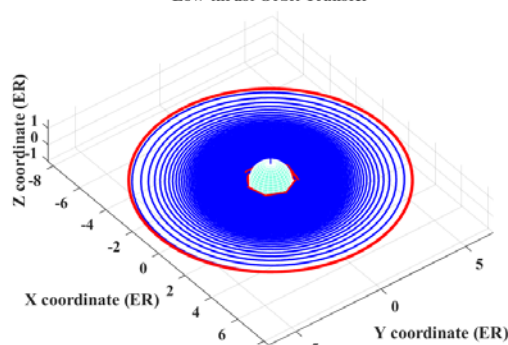
شکل ۲۰. تغییرات نیم قطر اصلی در حالت سوم مانور افزایش مدار

### ۵- نتیجه‌گیری

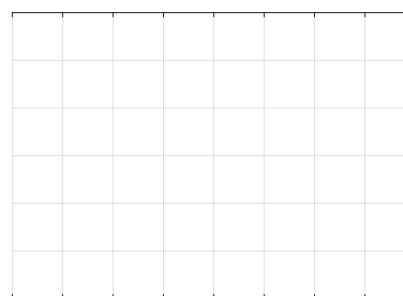
در مطالعات گذشته، استفاده از الگوریتم‌های یادگیری ماشین به خصوص یادگیری تقویتی در طراحی مسیر فضاپیمای تراست پایین در مأموریت‌های مانور مداری مانند ملاقات مداری، انتقال مداری و غیره مورد توجه قرار نگرفته است. در این مقاله، دینامیک مسئله طراحی مسیر فضاپیمای تراست پایین تحت قالب معادلات دیفرانسیل اعتدالی بر حسب عناصر مداری اعتدالی برای مدارهای دایروی هم صفحه بیان شد. با توجه به اینکه در مانور انتقال مداری هدف اصلی رسیدن به ارتفاع و خروج از مرکز مدار مورد نظر است، تابع پاداش به صورت مجموع وزن دهی شده عناصر مداری اعتدالی بیان شد. سپس الگوریتم شبکه بازیگر-منتقد برای سه حالت مانور افزایش مدار با دو نیروی تراست متفاوت به منظور نمایش عملکرد صحیح الگوریتم پیاده سازی شد. اثر تغییر فراپارامتر ضریب تنزل نیز بر روند یادگیری عامل به منظور تحلیل حساسیت الگوریتم نشان داده شد. حالت سوم افزایش ارتفاع از ۱۵۰۰ کیلومتری به مدار ژئو برای نمایش دقت و مقایسه نتایج به دست آمده در حالت‌های قبل پیاده‌سازی شد. در نهایت با بررسی نتایج به دست آمده از سه شبیه‌سازی انجام شده، عملکرد این نوع از الگوریتم‌های یادگیری ماشین در انجام مأموریت‌های انتقال مداری به خوبی نمایش داده شد. مقادیر مهم عناصر مداری مانند نیم قطر اصلی و خروج از مرکز مدارها که در مانور افزایش مدار لئو به ژئو اهمیت بالایی دارد در هر سه حالت در محدوده مقادیر صحیح تعیین شده و منطقی خود قرار دارند. در حالت اول عامل به دلیل عدم آگاهی از دینامیک مسئله نیاز به طی کردن اپیزودها و تعامل بیشتر با محیط دارد که این موضوع در نمودار پاداش مربوط به آن مشاهده و توصیف شد.



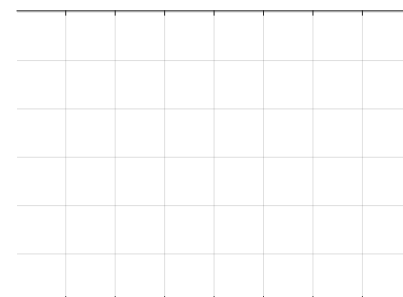
شکل ۱۶. نمودار پاداش در حالت سوم مانور افزایش ارتفاع مدار  
Low-thrust Orbit Transfer



شکل ۱۷. تغییرات نیم قطر اصلی در حالت دوم مانور افزایش مدار



شکل ۱۸. پروفیل زاویه انحراف تراستر در حالت سوم مانور افزایش مدار



شکل ۱۹. تغییرات خروج از مرکز در حالت سوم مانور افزایش مدار

نه تنها به نتایج دقت کافی را ندارد بلکه زمان زیادی نیز صرف یک مسیر اشتباه خواهد شد. همچنین در زمینه بهینه‌سازی، زمانی که اهداف مسئله از یک به دو یا بیشتر تبدیل می‌شود، تعریف و تنظیم کردن تابع پاداش و ضرایب آن و همچنین پارامترهای الگوریتم دشوار و پیچیده می‌شود که نیازمند تسلط کامل بر دینامیک مسئله و عملکرد الگوریتم در این دسته از مسائل است.

همچنین به عنوان یک کار مستقل و جذاب بعدی، می‌توان در کنار مسئله طراحی مسیر، پارامتر زمان و یا مصرف سوخت را نیز در معادلات و فرآیند پیاده‌سازی الگوریتم در نظر گرفت و با تنظیم تابع پاداش و شبکه‌های عصبی الگوریتم بازیگر منتقد این موضوع را نیز مورد بررسی و تحلیل قرار داد که به نوبه خود یک نیاز مهم در عملیات‌های واقعی مأموریت‌های مداری به حساب می‌آید.

### تعارض منافع

هیچ‌گونه تعارض منافع توسط نویسندگان بیان نشده است.

### مراجع

- [1] Robert E. Pritchett, "Numerical methods for low-thrust trajectory optimization", Purdue University West Lafayette, Indiana, August 2016.
- [2] J. B. Caillau, J. Gergaud, and J. Noailles. 3D geosynchronous transfer of a satellite: Continuation on the thrust. *Journal of Optimization Theory and Applications*, 118 (3):541–565, September 2003.
- [3] Mischa Kim, "Continuous Low-Thrust Trajectory Optimization: Techniques and Applications", Doctor of Philosophy Dissertation in Aerospace Engineering, April 18, 2005
- [4] Ryan P. Russell, "Primer Vector Theory Applied to Global Low-Thrust Trade Studies", California Institute of Technology, Pasadena, California 91109, *JOURNAL OF GUIDANCE, CONTROL, AND DYNAMICS*, 2007, DOI: 10.2514/1.22984
- [5] Weipeng Li Hai Huang, "Optimization of low-thrust transfers to libration point orbits", *Aircraft Engineering and Aerospace Technology: An International Journal* (2016), Vol. 88 Iss 1 pp. 16 – 23
- [6] Mauro Pontani · Bruce Conway, "Optimal Low-Thrust Orbital Maneuvers via Indirect Swarming Method" Accepted: 29 October 2013, Springer Science+Business Media New York 2013, DOI 10.1007/s10957-013-0471-9
- [7] P. Enright and B. A. Conway, "Discrete approximations to optimal trajectories using direct transcription and Nonlinear Programming," *Journal of Guidance Control and Dynamics*, Vol. 15, 09 1992, 10.2514/3.20934.
- [8] Jonathan D. Aziz, Jeffrey S. Parker, Daniel J. Scheeres, Jacob A. Englander, "Low-Thrust Many-Revolution Trajectory Optimization via Differential Dynamic Programming and a

ادامه حالت‌های دوم و سوم عامل با استفاده از تجربیات قبلی خود بهتر و سریعتر توانسته است به هدف تعیین شده با استفاده از تابع پاداش تعریف شده برسد که به خوبی در شکل‌های نمودار پاداش و عناصر مداری هر یک از این حالت‌ها قابل مشاهده است. از طرفی نیز می‌توان دید که با تغییر شرایط اولیه مانند نیروی تراست و ضرایب تنزل و یادگیری الگوریتم دقت این روش در مقایسه با روش‌های کنترل بهینه [۱۶] بیشتر است و حساسیت آن به این تغییرات نیز کمتر می‌باشد. همچنین روش‌های یادگیری ماشین جستجوهای جهانی را برای یافتن راه حل انجام می‌دهند که در مقایسه با سایر روش‌های بیان شده در مرور بر ادبیات یک مزیت محسوب می‌شود.

### جدول ۷. مقایسه دو رویکرد کنترل بهینه و یادگیری تقویتی [۱۷]

یادگیری تقویتی	کنترل بهینه
یک بهینه‌سازی تصادفی	یک بهینه‌سازی قطعی
زمان اجرای بسیار سریع برای سیاست آموزش دیده	زمان اجرا نامحدود به جز موارد خاص مانند محدودیت‌های تحذب
هیچ محدودیتی در نمایش پویایی مشکل وجود ندارد	دینامیک باید تحت یک معادله دیفرانسیل معمولی نشان داده شود.
حلقه بسته (هدایت و کنترل یکپارچه)	حلقه باز (برای ردیابی مسیر بهینه به یک کنترلر نیاز دارد)

مزیت اصلی و قابل توجه الگوریتم‌های یادگیری ماشین که باعث شده است تمرکز محققان به استفاده بیشتر از این روش‌ها معطوف شود این است که این نوع از الگوریتم‌ها بعد از آموزش دیدن بر روی دینامیک و محیط مسئله می‌توانند برای شرایط اولیه و قیود مختلفی که بر مأموریت حاکم می‌شود به کار گرفته شوند و دیگر نیاز به شبیه‌سازی و حل مجدد و کامل مسئله و تنظیم پارامترهای آن نیست که این خاصیت با بیان خودکار و خودران بودن این دسته از الگوریتم‌ها شناخته می‌شود.

از سوی دیگر، یکی از چالش‌های مهمی که روش یادگیری تقویتی با آن مواجه و بخش مهم الگوریتم است، تعیین تابع پاداش متناسب با دینامیک مسئله مربوطه است. چراکه زمانی که تابع پاداش خیلی پراکنده و دور از هدف اصلی مسئله باشد

[19] Navaee, M., & Sanati, M. (2013). Optimal Impulsive Orbital 3D Maneuver with or without Time Constraint. *Amirkabir Journal of Mechanical Engineering*, 44(2), 53-69.

[20] Fakoor, M., Sadeghi, S., & Bakhtiari, M. (2020). Investigation of low thrust optimal orbital transfer from LEO to GEO considering circular orbits. *The Journal of the Astronautical Sciences*, 67, 77-97.



## COPYRIGHTS

© 2024 by the authors. Licensee Iranian Space Research Center of Iran. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution 4.0 International (CC BY 4.0) (<https://creativecommons.org/licenses/by/4.0/>)

Sundman Transformation”, American Astronautical Society, part of Springer Nature 2018

[9] S. N. Williams and V. Coverstone-Carroll. Mars missions using solar electric propulsion. *Journal of Spacecraft and Rockets*, 37(1):71–77, January–February 2000.

[10] Zhenbo Wang and Michael J. Grant, “Minimum-Fuel Low-Thrust Transfers for Spacecraft: A Convex Approach” School of Aeronautics and Astronautics, Purdue University, West Lafayette, Indiana, 47907-2045, USA, DOI 10.1109/TAES.2018.2812558, IEEE Transactions on Aerospace and Electronic Systems 2018

[11] Daniel S. Kolosa, “A Reinforcement Learning Approach to Spacecraft Trajectory Optimization”, Western Michigan University, Dissertations. 3542, 2019

[12] M.P. van Hoom, “OPTIMIZING AIR-TO-AIR MISSILE GUIDANCE USING REINFORCEMENT LEARNING”, Master of Science thesis in Aerospace Engineering, Delft University of Technology March 26, 2019

[13] T.A.H. Kranen, “Low-Thrust Gravity Assist Trajectory Optimisation using Evolutionary Neurocontrol”, Dissertation TU Delft Aerospace Engineering, 2019.

[14] Howard D. Curtis, “*Orbital Mechanics for Engineering Students*”, Third Edition, season 6, p.g 344, 2014

[15] Richard S. Sutton and Andrew G. Barto, Reinforcement Learning: An Introduction, MIT Press, Cambridge, MA, 1998A Bradford Book

[16] M. Navabi, E. Meshkinfam, "Space low-thrust trajectory optimization utilizing numerical techniques, a comparative study," 2013 6th International Conference on Recent Advances in Space Technologies (RAST), IEEE, 2013, pp. 303-307, doi: 10.1109/RAST.2013.6581222.

[17] Brian Gaudet, Roberto Furfaro, Richard Linares, “Reinforcement Learning for Angle-Only Intercept Guidance of Maneuvering Targets”, *Aerospace Science and Technology*, Volume 99, April 2020

[18] Navabi, M., & Sabatifar, M. (2010). Optimal Impulsive Maneuver Between Elliptical Coplanar-Noncoplanar Orbits. *Space Science and Technology*, 3(1), 67-74.

xvii Agent

xviii Policy

xix reward value

xx Value

xxi expected return

xxii discount factor

xxiii policy gradients

xxiv Deep Deterministic Policy Gradients (DDPG)

xxv state-value function

xxvi temporal difference error

xxvii Gauss’s Variational equations

xxviii equinoctial orbital elements

xxix model-free

xxx Exploration

xxxi Exploitation

i Orbital Rendezvous

ii orbit maintenance

iii specific impulse(Isp)

iv Two-boundary value problem

v Pontryagin

vi Adaptive Simulated Annealing

vii Particle Swarm Optimization

viii Sundman

ix differential dynamic programming

x Genetic Algorithm

xi Artificial Bee Colony (ABC) algorithm

xii initial value problems

xiii Reinforcement Learning(RL)

xiv Actor-Critic networks

xv Supervised Learning (SL)

xvi Unsupervised Learning (UL)